



Feynman's Preface

These are the lectures in physics that I gave last year and the year before to the freshman and sophomore classes at Caltech. The lectures are, of course, not verbatim—they have been edited, sometimes extensively and sometimes less so. The lectures form only part of the complete course. The whole group of 180 students gathered in a big lecture room twice a week to hear these lectures and then they broke up into small groups of 15 to 20 students in recitation sections under the guidance of a teaching assistant. In addition, there was a laboratory session once a week.

The special problem we tried to get at with these lectures was to maintain the interest of the very enthusiastic and rather smart students coming out of the high schools and into Caltech. They have heard a lot about how interesting and exciting physics is—the theory of relativity, quantum mechanics, and other modern ideas. By the end of two years of our previous course, many would be very discouraged because there were really very few grand, new, modern ideas presented to them. They were made to study inclined planes, electrostatics, and so forth, and after two years it was quite stultifying. The problem was whether or not we could make a course which would save the more advanced and excited student by maintaining his enthusiasm.

The lectures here are not in any way meant to be a survey course, but are very serious. I thought to address them to the most intelligent in the class and to make sure, if possible, that even the most intelligent student was unable to completely encompass everything that was in the lectures—by putting in suggestions of applications of the ideas and concepts in various directions outside the main line of attack. For this reason, though, I tried very hard to make all the statements as accurate as possible, to point out in every case where the equations and ideas fitted into the body of physics, and how—when they learned more—things would be modified. I also felt that for such students it is important to indicate what it is that they should—if they are sufficiently clever—be able to understand by deduction from what has been said before, and what is being put in as something new. When new ideas came in, I would try either to deduce them if they were deducible, or to explain that it *was* a new idea which hadn't any basis in terms of things they had already learned and which was not supposed to be provable—but was just added in.

At the start of these lectures, I assumed that the students knew something when they came out of high school—such things as geometrical optics, simple chemistry ideas, and so on. I also didn't see that there was any reason to make the lectures

in a definite order, in the sense that I would not be allowed to mention something until I was ready to discuss it in detail. There was a great deal of mention of things **to** come, without complete discussions. These more complete discussions would come later when the preparation became more advanced. Examples are the discussions of inductance, and of energy levels, which are at first brought in in a very qualitative way and are later developed more completely.

At the same time that I was aiming at the more active student, I also wanted to take care of the fellow for whom the extra fireworks and side applications are merely disquieting and who cannot be expected to learn most of the material in the lecture at all. For such students I wanted there to be at least a central core or backbone of material which he *could* get. Even if he didn't understand everything in a lecture, I hoped he wouldn't get nervous. I didn't expect him to understand everything, but only the central and most direct features. It takes, of course, a certain intelligence on his part to see which are the central theorems and central ideas, and which are the more advanced side issues and applications which he may understand only in later years.

In giving these lectures there was one serious difficulty: in the way the course was given, there wasn't any feedback from the students to the lecturer to indicate how well the lectures were going over. This is indeed a very serious difficulty, and I don't know how good the lectures really are. The whole thing was essentially an experiment. And if I did it again I wouldn't do it the same way—I hope I *don't* have to do it again! I think, though, that things worked out—so far as the physics is concerned—quite satisfactorily in the first year.

In the second year I was not so satisfied. In the first part of the course, dealing with electricity and magnetism, I couldn't think of any really unique or different way of doing it—of any way that would be particularly more exciting than the usual way of presenting it. So I don't think I did very much in the lectures on electricity and magnetism. At the end of the second year I had originally intended to go on, after the electricity and magnetism, by giving some more lectures on the properties of materials, but mainly to take up things like fundamental modes, solutions of the diffusion equation, vibrating systems, orthogonal functions,... developing the first stages of what are usually called "the mathematical methods of physics." In retrospect, I think that if I were doing it again I would go back to that original idea. But since it was not planned that I would be giving these lectures again, it was suggested that it might be a good idea to try to give an introduction to the quantum mechanics—what you will find in Volume III.

It is perfectly clear that students who will major in physics can wait until their third year for quantum mechanics. On the other hand, the argument was made that many of the students in our course study physics as a background for their primary interest in other fields. And the usual way of dealing with quantum mechanics makes that subject almost unavailable for the great majority of students because they have to take so long to learn it. Yet, in its real applications—especially in its more complex applications, such as in electrical engineering and chemistry—the full machinery of the differential equation approach is not actually used. So I tried to describe the principles of quantum mechanics in a way which wouldn't require that one first know the mathematics of partial differential equations. Even for a physicist I think that is an interesting thing to try to do—to present quantum mechanics in this reverse fashion—for several reasons which may be apparent in the lectures themselves. However, I think that the experiment in the quantum mechanics part was not completely successful—in large part because I really did not have enough time at the end (I should, for instance, have had three or four more lectures in order to deal more completely with such matters as energy bands and the spatial dependence of amplitudes). Also, I had never presented the subject this way before, so the lack of feedback was particularly serious. I now believe the quantum mechanics should be given at a later time. Maybe I'll have a chance to do it again someday. Then I'll do it right.

The reason there are no lectures on how to solve problems is because there **were** recitation sections. Although I did put in three lectures in the first year on how to solve problems, they are not included here. Also there was a lecture on inertial

guidance which certainly belongs after the lecture on rotating systems, but which was, unfortunately, omitted. The fifth and sixth lectures are actually due to Matthew Sands, as I was out of town.

The question, of course, is how well this experiment has succeeded. My own point of view—which, however, does not seem to be shared by most of the people who worked with the students—is pessimistic. I don't think I did very well by the students. When I look at the way the majority of the students handled the problems on the examinations, I think that the system is a failure. Of course, my friends point out to me that there were one or two dozen students who—very surprisingly—understood almost everything in all of the lectures, and who were quite active in working with the material and worrying about the many points in an excited and interested way. These people have now, I believe, a first-rate background in physics—and they are, after all, the ones I was trying to get at. But then, "The power of instruction is seldom of much efficacy except in those happy dispositions where it is almost superfluous." (Gibbon)

Still, I didn't want to leave any student completely behind, as perhaps I did. I think one way we could help the students more would be by putting more hard work into developing a set of problems which would elucidate some of the ideas in the lectures. Problems give a good opportunity to fill out the material of the lectures and make more realistic, more complete, and more settled in the mind the ideas that have been exposed.

I think, however, that there isn't any solution to this problem of education other than to realize that the best teaching can be done only when there is a direct individual relationship between a student and a good teacher—a situation in which the student discusses the ideas, thinks about the things, and talks about the things. It's impossible to learn very much by simply sitting in a lecture, or even by simply doing problems that are assigned. But in our modern times we have so many students to teach that we have to try to find some substitute for the ideal. Perhaps my lectures can make some contribution. Perhaps in some small place where there are individual teachers and students, they may get some inspiration or some ideas from the lectures. Perhaps they will have fun thinking them through—or going on to develop some of the ideas further.

RICHARD P. FEYNMAN

June, 1963

Foreword

This book is based upon a course of lectures in introductory physics given by Prof. R. P. Feynman at the California Institute of Technology during the academic year 1961-62; it covers the first year of the two-year introductory course taken by all Caltech freshmen and sophomores, and was followed in 1962-63 by a similar series covering the second year. The lectures constitute a major part of a fundamental revision of the introductory course, carried out over a four-year period.

The need for a basic revision arose both from the rapid development of physics in recent decades and from the fact that entering freshmen have shown a steady increase in mathematical ability as a result of improvements in high school mathematics course content. We hoped to take advantage of this improved mathematical background, and also to introduce enough modern subject matter to make the course challenging, interesting, and more representative of present-day physics.

In order to generate a variety of ideas on what material to include and how to present it, a substantial number of the physics faculty were encouraged to offer their ideas in the form of topical outlines for a revised course. Several of these were presented and were thoroughly and critically discussed. It was agreed almost at once that a basic revision of the course could not be accomplished either by merely adopting a different textbook, or even by writing one *ab initio*, but that the new course should be centered about a set of lectures, to be presented at the rate of two or three per week; the appropriate text material would then be produced as a secondary operation as the course developed, and suitable laboratory experiments would also be arranged to fit the lecture material. Accordingly, a rough outline of the course was established, but this was recognized as being incomplete, tentative, and subject to considerable modification by whoever was to bear the responsibility for actually preparing the lectures.

Concerning the mechanism by which the course would finally be brought to life, several plans were considered. These plans were mostly rather similar, involving a cooperative effort by N staff members who would share the total burden symmetrically and equally: each man would take responsibility for $1/N$ of the material, deliver the lectures, and write text material for his part. However, the unavailability of sufficient staff, and the difficulty of maintaining a uniform point of view because of differences in personality and philosophy of individual participants, made such plans seem unworkable.

The realization that we actually possessed the means to create not just a new and different physics course, but possibly a unique one, came as a happy inspiration to Professor Sands. He suggested that Professor R. P. Feynman prepare and deliver the lectures, and that these be tape-recorded. When transcribed and edited, they would then become the textbook for the new course. This is essentially the plan that was adopted.

It was expected that the necessary editing would be minor, mainly consisting of supplying figures, and checking punctuation and grammar; it was to be done by one or two graduate students on a part-time basis. Unfortunately, this expectation was short-lived. It was, in fact, a major editorial operation to transform the verbatim transcript into readable form, even without the reorganization or revision of the subject matter that was sometimes required. Furthermore, it was not a job for a technical editor or for a graduate student, but one that required the close attention of a professional physicist for from ten to twenty hours per lecture!

The difficulty of the editorial task, together with the need to place the material in the hands of the students as soon as possible, set a strict limit upon the amount of "polishing" of the material that could be accomplished, and thus we were forced to aim toward a preliminary but technically correct product that could be used immediately, rather than one that might be considered final or finished. Because of an urgent need for more copies for our students, and a heartening interest on the part of instructors and students at several other institutions, we decided to publish the material in its preliminary form rather than wait for a further major revision which might never occur. We have no illusions as to the completeness, smoothness, or logical organization of the material; in fact, we plan several minor modifications in the course in the immediate future, and we hope that it will not become static in form or content.

In addition to the lectures, which constitute a centrally important part of the course, it was necessary also to provide suitable exercises to develop the students' experience and ability, and suitable experiments to provide first-hand contact with the lecture material in the laboratory. Neither of these aspects is in as advanced a state as the lecture material, but considerable progress has been made. Some exercises were made up as the lectures progressed, and these were expanded and amplified for use in the following year. However, because we are not yet satisfied that the exercises provide sufficient variety and depth of application of the lecture material to make the student fully aware of the tremendous power being placed at his disposal, the exercises are published separately in a less permanent form in order to encourage frequent revision.

A number of new experiments for the new course have been devised by Professor H. V. Neher. Among these are several which utilize the extremely low friction exhibited by a gas bearing: a novel linear air trough, with which quantitative measurements of one-dimensional motion, impacts, and harmonic motion can be made, and an air-supported, air-driven Maxwell top, with which accelerated rotational motion and gyroscopic precession and nutation can be studied. The development of new laboratory experiments is expected to continue **for a** considerable period of time.

The revision program was under the direction of Professors R. B. Leighton, H. V. Neher, and M. Sands. Officially participating in the program were Professors R. P. Feynman, G. Neugebauer, R. M. Sutton, H. P. Stabler,* F. Strong, and R. Vogt, from the division of Physics, Mathematics and Astronomy, and Professors T. Caughey, M. Plesset, and C. H. Wilts from the division of Engineering Science. The valuable assistance of all those contributing to the revision program is gratefully acknowledged. We are particularly indebted to the Ford Foundation, without whose financial assistance this program could not have been carried out.

ROBERT B. LEIGHTON

July, 1963

* 1961-62, while on leave from Williams College, Williamstown, Mass.

Contents

CHAPTER 1. ATOMS IN MOTION

- 1-1 Introduction 1-1
- 1-2 Matter is made of atoms 1-2
- 1-3 Atomic processes 1-5
- 1-4 Chemical reactions 1-6

CHAPTER 2. BASIC PHYSICS

- 2-1 Introduction 2-1
- 2-2 Physics before 1920 2-3
- 2-3 Quantum physics 2-6
- 2-4 Nuclei and particles 2-8

CHAPTER 3. THE RELATION OF PHYSICS TO OTHER SCIENCES

- 3-1 Introduction 3-1
- 3-2 Chemistry 3-1
- 3-3 Biology 3-2
- 3-4 Astronomy 3-6
- 3-5 Geology 3-7
- 3-6 Psychology 3-8
- 3-7 How did it get that way? 3-9

CHAPTER 4. CONSERVATION OF ENERGY

- 4-1 What is energy? 4-1
- 4-2 Gravitational potential energy 4-2
- 4-3 Kinetic energy 4-5
- 4-4 Other forms of energy 4-6**

CHAPTER 5. TIME AND DISTANCE

- 5-1 Motion 5-1
- 5-2 Time 5-1
- 5-3 Short times 5-2
- 5-4 Long times 5-3
- 5-5 Units and standards of time 5-5
- 5-6 Large distances 5-5
- 5-7 Short distances 5-8

CHAPTER 6. PROBABILITY

- 6-1 Chance and likelihood 6-1
- 6-2 Fluctuations 6-3
- 6-3 The random walk 6-5
- 6-4 A probability distribution 6-7
- 6-5 The uncertainty principle 6-10

CHAPTER 7. THE THEORY OF GRAVITATION

- 7-1 Planetary motions 7-1
- 7-2 Kepler's laws 7-1
- 7-3 Development of dynamics 7-2
- 7-4 Newton's law of gravitation 7-3**
- 7-5 Universal gravitation 7-5
- 7-6 Cavendish's experiment 7-9
- 7-7 What is gravity? 7-9
- 7-8 Gravity and relativity 7-11

CHAPTER 8. MOTION

- 8-1 Description of motion 8-1
- 8-2 Speed 8-2
- 8-3 Speed as a derivative 8-5
- 8-4 Distance as an integral 8-7
- 8-5 Acceleration 8-8

CHAPTER 9. NEWTON'S LAWS OF DYNAMICS

- 9-1 Momentum and force 9-1
- 9-2 Speed and velocity 9-2
- 9-3 Components of velocity, acceleration, and force 9-3**
- 9-4 What is the force? 9-3**
- 9-5 Meaning of the dynamical equations 9-4
- 9-6 Numerical solution of the equations 9-5
- 9-7 Planetary motions 9-6

CHAPTER 10. CONSERVATION OF MOMENTUM

- 10-1 Newton's Third Law 10-1
- 10-2 Conservation of momentum 10-2
- 10-3 Momentum is conserved! 10-5
- 10-4 Momentum and energy 10-7
- 10-5 Relativistic momentum 10-8

CHAPTER 11. VECTORS

- 11-1 Symmetry in physics **11-1**
- 11-2 Translations 11-1
- 11-3 Rotations 11-3
- 11-4 Vectors 11-5**
- 11-5 Vector algebra 11-6**
- 11-6 Newton's laws in vector notation **11-7**
- 11-7 Scalar product of vectors 11-8

CHAPTER 12. CHARACTERISTICS OF FORCE

- 12-1 What is a force? 12-1
- 12-2 Friction 12-3
- 12-3 Molecular forces 12-6
- 12-4 Fundamental forces. Fields 12-7
- 12-5 Pseudo forces 12-10
- 12-6 Nuclear forces 12-12

CHAPTER 13. WORK AND POTENTIAL ENERGY (A)

- 13-1 Energy of a falling body 13-1
- 13-2 Work done by gravity 13-3
- 13-3 Summation of energy 13-6
- 13-4 Gravitational field of large objects 13-8

CHAPTER 14. WORK AND POTENTIAL ENERGY (conclusion)

- 14-1 Work 14-1
- 14-2 Constrained motion 14-3
- 14-3 Conservative forces 14-3
- 14-4 Nonconservative forces 14-6
- 14-5 Potentials and fields 14-7

CHAPTER 15. THE SPECIAL THEORY OF RELATIVITY

- 15-1** The principle of relativity 15-1
- 15-2** The Lorentz transformation 15-3
- 15-3** The Michelson-Morley experiment 15-3
- 15-4** Transformation of time 15-5
- 15-5** The Lorentz contraction 15-7
- 15-6** Simultaneity 15-7
- 15-7 Four-vectors 15-8
- 15-8 Relativistic dynamics 15-9
- 15-9** Equivalence of mass and energy 15-10

CHAPTER 16. RELATIVISTIC ENERGY AND MOMENTUM

- 16-1** Relativity and the philosophers **16-1**
- 16-2 The twin paradox 16-3
- 16-3 Transformation of velocities 16-4
- 16-4 Relativistic mass 16-6
- 16-5 Relativistic energy 16-8

CHAPTER 17. SPACE-TIME

- 17-1 The geometry of space-time 17-1
- 17-2 Space-time intervals 17-2
- 17-3 Past, present, and future 17-4
- 17-4 More about four-vectors 17-5
- 17-5 Four-vector algebra 17-7.

CHAPTER 18. ROTATION IN Two DIMENSIONS

- 18-1 The center of mass 18-1
- 18-2 Rotation of a rigid body 18-2
- 18-3 Angular momentum 18-5
- 18-4 Conservation of angular momentum 18-6

CHAPTER 19. CENTER OF MASS; MOMENT OF INERTIA

- 19-1 Properties of the center of mass 19-1
- 19-2 Locating the center of mass 19-4
- 19-3 Finding the moment of inertia 19-5
- 19-4** Rotational kinetic energy 19-7

CHAPTER 20. ROTATION IN SPACE

- 20-1 Torques in three dimensions 20-1
- 20-2 The rotation equations using cross **products 20-4**
- 20-3 The gyroscope 20-5
- 20-4 Angular momentum of a solid body 20-8

CHAPTER 21. THE HARMONIC OSCILLATOR

- 21-1 Linear differential equations 21-1
- 21-2 The harmonic oscillator 21-1
- 21-3 Harmonic motion and circular motion 21-4
- 21-4 Initial conditions 21-4
- 21-5 Forced oscillations 21-5

CHAPTER 22. ALGEBRA

- 22-1 Addition and multiplication 22-1
- 22-2 The inverse operations 22-2
- 22-3 Abstraction and generalization 22-3
- 22-4 Approximating irrational numbers 22-4
- 22-5 Complex numbers 22-7
- 22-6 Imaginary exponents 22-9

CHAPTER 23. RESONANCE

- 23-1 Complex numbers and harmonic motion 23-1
- 23-2 The forced oscillator with damping 23-3

- 23-3 Electrical resonance 23-5
- 23-4 Resonance in nature 23-7

CHAPTER 24. TRANSIENTS

- 24-1 The energy of an oscillator 24-1
- 24-2 Damped oscillations 24-2
- 24-3 Electrical transients 24-5

CHAPTER 25. LINEAR SYSTEMS AND REVIEW

- 25-1 Linear differential equations 25-1
- 25-2 Superposition of solutions 25-2
- 25-3 Oscillations in linear systems 25-5
- 25-4 Analogs in physics 25-6
- 25-5 Series and parallel impedances 25-8

CHAPTER 26. OPTICS: THE PRINCIPLE OF LEAST TIME

- 26-1 Light 26-1
- 26-2 Reflection and refraction 26-2
- 26-3 Fermat's principle of least time 26-3
- 26-4 Applications of Fermat's principle 26-5
- 26-5 A more precise statement of Fermat's principle 26-7
- 26-6 How it works 26-8

CHAPTER 27. GEOMETRICAL OPTICS

- 27-1 Introduction 27-1
- 27-2 The focal length of a spherical surface 27-1
- 27-3 The focal length of a lens 27-4
- 27-4 Magnification 27-5
- 27-5 Compound lenses **27-6**
- 27-6 Aberrations 27-7
- 27-7 Resolving power 27-7

CHAPTER 28. ELECTROMAGNETIC RADIATION

- 28-1 Electromagnetism 28-1
- 28-2 Radiation 28-3
- 28-3 **The** dipole radiator 28-5
- 28-4 Interference 28-6

CHAPTER 29. INTERFERENCE

- 29-1 Electromagnetic waves 29-1
- 29-2 Energy of radiation 29-2
- 29-3 Sinusoidal waves 29-2
- 29-4 Two dipole radiators 29-3
- 29-5 The mathematics of interference 29-5

CHAPTER 30. DIFFRACTION

- 30-1 The resultant amplitude due to n equal oscillators 30-1
- 30-2 The diffraction grating 30-3
- 30-3 Resolving power of a grating **30-5**
- 30-4 The parabolic antenna 30-6
- 30-5 Colored films; crystals 30-7
- 30-6 Diffraction by opaque screens 30-8
- 30-7 The field of a plane of oscillating charges 30-10

CHAPTER 31. THE ORIGIN OF THE REFRACTIVE INDEX

- 31-1 The index of refraction 31-1
- 31-2 The field due to the material 31-4
- 31-3 Dispersion 31-6
- 31-4 Absorption 31-8
- 31-5 The energy carried by an electric wave 31-9
- 31-6 Diffraction of light by a screen 31-10

CHAPTER 32. RADIATION DAMPING. **LIGHT** SCATTERING

- 32-1 Radiation resistance 32-1
- 32-2 The rate of radiation of energy 32-2
- 32-3 Radiation damping 32-3
- 32-4 Independent sources 32-5
- 32-5 Scattering of light 32-6

CHAPTER 33. POLARIZATION

- 33-1 The electric vector of light 33-1
- 33-2 Polarization of scattered light 33-3
- 33-3 Birefringence 33-3
- 33-4 Polarizers 33-5
- 33-5 Optical activity 33-6
- 33-6 The intensity of reflected light 33-7
- 33-7 Anomalous refraction 33-9

CHAPTER 34. RELATIVISTIC EFFECTS IN RADIATION

- 34-1 Moving sources 34-1
- 34-2 Finding the "apparent" motion 34-2
- 34-3 Synchrotron radiation 34-3
- 34-4 Cosmic synchrotron radiation 34-6
- 34-5 Bremsstrahlung 34-6
- 34-6 The Doppler effect 34-7
- 34-7 The ω , k four-vector 34-9
- 34-8 Aberration 34-10
- 34-9 The momentum of light 34-10

CHAPTER 35. COLOR VISION

- 35-1 The human eye 35-1
- 35-2 Color depends on intensity 35-2
- 35-3 Measuring the color sensation 35-3
- 35-4 The chromaticity diagram 35-6
- 35-5 The mechanism of color vision 35-7
- 35-6 Physiochemistry of color vision 35-9

CHAPTER 36. MECHANISMS OF SEEING

- 36-1 The sensation of color 36-1
- 36-2 The physiology of the eye 36-3
- 36-3 The rod cells 36-6
- 36-4 The compound (insect) eye 36-6
- 36-5 Other eyes 36-9
- 36-6 Neurology of vision 36-9

CHAPTER 37. QUANTUM BEHAVIOR

- 37-1 Atomic mechanics 37-1
- 37-2 An experiment with bullets 37-2
- 37-3 An experiment with waves 37-3
- 37-4 An experiment with electrons 37-4
- 37-5 The interference of electron waves 37-5
- 37-6 Watching the electrons 37-7
- 37-7 First principles of quantum mechanics 37-10
- 37-8 The uncertainty principle 37-11

CHAPTER 38. THE RELATION OF WAVE AND PARTICLE VIEWPOINTS

- 38-1 Probability wave amplitudes 38-1
- 38-2 Measurement of position and momentum 38-2
- 38-3 Crystal diffraction 38-4
- 38-4 The size of an atom 38-5

- 38-5 Energy levels 38-7
- 38-6 Philosophical implications 38-8

CHAPTER 39. THE KINETIC THEORY OF GASES

- 39-1 Properties of matter 39-1
- 39-2 The pressure of a gas 39-2
- 39-3 Compressibility of radiation 39-6
- 39-4 Temperature and kinetic energy 39-6
- 39-5 The ideal gas law 39-10

CHAPTER 40. THE PRINCIPLES OF STATISTICAL MECHANICS

- 40-1 The exponential atmosphere 40-1
- 40-2 The Boltzmann law 40-2
- 40-3 Evaporation of a liquid 40-3
- 40-4 The distribution of molecular speeds 40-4
- 40-5 The specific heats of gases 40-7
- 40-6 The failure of classical physics 40-8

CHAPTER 41. THE BROWNIAN MOVEMENT

- 41-1 Equipartition of energy 41-1
- 41-2 Thermal equilibrium of radiation 41-3
- 41-3 Equipartition and the quantum oscillator 41-6
- 41-4** The random walk 41-8

CHAPTER 42. APPLICATIONS OF KINETIC THEORY

- 42-1 Evaporation 42-1
- 42-2 Thermionic emission 42-4
- 42-3 Thermal ionization 42-5
- 42-4** Chemical kinetics 42-7
- 42-5 Einstein's laws of radiation 42-8

CHAPTER 43. DIFFUSION

- 43-1 Collisions between molecules 43-1
- 43-2 The mean free path 43-3
- 43-3 The drift speed 43-4
- 43-4 Ionic conductivity 43-6
- 43-5 Molecular diffusion 43-7
- 43-6 Thermal conductivity 43-9

CHAPTER 44. THE LAWS OF THERMODYNAMICS

- 44-1 Heat engines; the first law 44-1
- 44-2 The second law 44-3
- 44-3 Reversible engines 44-4
- 44-4 The efficiency of an ideal engine 44-7
- 44-5 The thermodynamic temperature 44-9
- 44-6 Entropy 44-10

CHAPTER 45. ILLUSTRATIONS OF THERMODYNAMICS

- 45-1 Internal energy 45-1
- 45-2 Applications 45-4
- 45-3 The Clausius-Clapeyron equation 45-6

CHAPTER 46. RATCHET AND PAWL

- 46-1 How a ratchet works 46-1
- 46-2 The ratchet as an engine 46-2
- 46-3 Reversibility in mechanics 46-4
- 46-4 Irreversibility 46-5
- 46-5 Order and entropy 46-7

CHAPTER 47. SOUND. THE WAVE EQUATION

- 47-1 Waves 47-1
- 47-2 The propagation of sound 47-3
- 47-3 The wave equation 47-4
- 47-4 Solutions of the wave equation 47-6
- 47-5 **The speed of sound** 47-7

CHAPTER 48. BEATS

- 48-1 Adding two waves 48-1
- 48-2 Beat notes and modulation 48-3
- 48-3 Side bands 48-4
- 48-4 Localized wave trains 48-5
- 48-5 Probability amplitudes for particles 48-7
- 48-6 Waves in three dimensions 48-9
- 48-7 Normal modes 48-10

CHAPTER 49. MODES

- 49-1 The reflection of waves 49-1
- 49-2 Confined waves, with natural frequencies 49-2
- 49-3 Modes in two dimensions 49-3
- 49-4 Coupled pendulums 49-6
- 49-5 Linear systems 49-7

INDEX

CHAPTER 50. HARMONICS

- 50-1 Musical tones 50-1
- 50-2 The Fourier series 50-2
- 50-3 Quality and consonance 50-3
- 50-4 The Fourier coefficients 50-5
- 50-5 The energy theorem 50-7
- 50-6 Nonlinear responses 50-8

CHAPTER 51. WAVES

- 51-1 Bow waves **51-1**
- 51-2 Shock waves 51-2
- 51-3 Waves in solids 51-4
- 51-4 Surface waves 51-7

CHAPTER 52. SYMMETRY IN PHYSICAL LAWS

- 52-1 Symmetry operations 52-1
- 52-2 Symmetry **in space and time** **52-1**
- 52-3 Symmetry and conservation laws 52-3
- 52-4 Mirror reflections 52-4
- 52-5 Polar and axial vectors 52-6
- 52-6 Which hand is right? 52-8
- 52-7 Parity is not conserved! 52-8
- 52-8 Antimatter 52-10
- 52-9 Broken symmetries 52-11

Atoms in Motion

1-1 Introduction

This two-year course in physics is presented from the point of view that you, the reader, are going to be a physicist. This is not necessarily the case of course, but that is what every professor in every subject assumes! If you are going to be a physicist, you will have a lot to study: two hundred years of the most rapidly developing field of knowledge that there is. So much knowledge, in fact, that you might think that you cannot learn all of it in four years, and truly you cannot; you will have to go to graduate school too!

Surprisingly enough, in spite of the tremendous amount of work that has been done for all this time it is possible to condense the enormous mass of results to a large extent—that is, to find *laws* which summarize all our knowledge. Even so, the laws are so hard to grasp that it is unfair to you to start exploring this tremendous subject without some kind of map or outline of the relationship of one part of the subject of science to another. Following these preliminary remarks, the first three chapters will therefore outline the relation of physics to the rest of the sciences, the relations of the sciences to each other, and the meaning of science, to help us develop a "feel" for the subject.

You might ask why we cannot teach physics by just giving the basic laws on page one and then showing how they work in all possible circumstances, as we do in Euclidean geometry, where we state the axioms and then make all sorts of deductions. (So, not satisfied to learn physics in four years, you want to learn it in four minutes?) We cannot do it in this way for two reasons. First, we do not yet *know* all the basic laws: there is an expanding frontier of ignorance. Second, the correct statement of the laws of physics involves some very unfamiliar ideas which require advanced mathematics for their description. Therefore, one needs a considerable amount of preparatory training even to learn what the *words* mean. No, it is not possible to do it that way. We can only do it piece by piece.

Each piece, or part, of the whole of nature is always merely an *approximation* to the complete truth, or the complete truth so far as we know it. In fact, everything we know is only some kind of approximation, because *we know that we do not know all the laws* as yet. Therefore, things must be learned only to be unlearned again or, more likely, to be corrected.

The principle of science, the definition, almost, is the following: *The test of all knowledge is experiment.* Experiment is the *sole judge* of scientific "truth." But what is the source of knowledge? Where do the laws that are to be tested come from? Experiment, itself, helps to produce these laws, in the sense that it gives us hints. But also needed is *imagination* to create from these hints the great generalizations—to guess at the wonderful, simple, but very strange patterns beneath them all, and then to experiment to check again whether we have made the right guess. This imagining process is so difficult that there is a division of labor in physics: there are *theoretical* physicists who imagine, deduce, and guess at new laws, but do not experiment; and then there are *experimental* physicists who experiment, imagine, deduce, and guess.

We said that the laws of nature are approximate: that we first find the "wrong" ones, and then we find the "right" ones. Now, how can an experiment be "wrong"? First, in a trivial way: if something is wrong with the apparatus that you did not notice. But these things are easily fixed, and checked back and forth. So without snatching at such minor things, how *can* the results of an experiment be wrong? Only by being inaccurate. For example, the mass of an object never seems to

1-1 Introduction

1-2 Matter is made of atoms

1-3 Atomic processes

1-4 Chemical reactions

change; a spinning top has the same weight as a still one. So a "law" was invented: mass is constant, independent of speed. That "law" is now found to be incorrect. Mass is found to increase with velocity, but appreciable increases require velocities near that of light. A *true* law is: if an object moves with a speed of less than one hundred miles a second the mass is constant to within one part in a million. In some such approximate form this is a correct law. So in practice one might think that the new law makes no significant difference. Well, yes and no. For ordinary speeds we can certainly forget it and use the simple constant-mass law as a good approximation. But for high speeds we are wrong, and the higher the speed, the more wrong we are.

Finally, and most interesting, *philosophically we are completely wrong* with the approximate law. Our entire picture of the world has to be altered even though the mass changes only by a little bit. This is a very peculiar thing about the philosophy, or the ideas, behind the laws. Even a very small effect sometimes requires profound changes in our ideas.

Now, what should we teach first? Should we teach the *correct* but unfamiliar law with its strange and difficult conceptual ideas, for example the theory of relativity, four-dimensional space-time, and so on? Or should we first teach the simple "constant-mass" law, which is only approximate, but does not involve such difficult ideas? The first is more exciting, more wonderful, and more fun, but the second is easier to get at first, and is a first step to a real understanding of the second idea. This point arises again and again in teaching physics. At different times we shall have to resolve it in different ways, but at each stage it is worth learning what is now known, how accurate it is, how it fits into everything else, and how it may be changed when we learn more.

Let us now proceed with our outline, or general map, of our understanding of science today (in particular, physics, but also of other sciences on the periphery), so that when we later concentrate on some particular point we will have some idea of the background, why that particular point is interesting, and how it fits into the big structure. So, what *is* our over-all picture of the world?

1-2 Matter is made of atoms

If, in some cataclysm, all of scientific knowledge were to be destroyed, and only one sentence passed on to the next generations of creatures, what statement would contain the most information in the fewest words? I believe it is the *atomic hypothesis* (or the *atomic fact*, or whatever you wish to call it) that *all things are made of atoms—little particles that move around in perpetual motion, attracting each other when they are a little distance apart, but repelling upon being squeezed into one another*. In that one sentence, you will see, there is an *enormous* amount of information about the world, if just a little imagination and thinking are applied.

To illustrate the power of the atomic idea, suppose that we have a drop of water a quarter of an inch on the side. If we look at it very closely we see nothing but water—smooth, continuous water. Even if we magnify it with the best optical microscope available—roughly two thousand times—then the water drop will be roughly forty feet across, about as big as a large room, and if we looked rather closely, we would *still* see relatively smooth water—but here and there small football-shaped things swimming back and forth. Very interesting. These are paramecia. You may stop at this point and get so curious about the paramecia with their wiggling cilia and twisting bodies that you go no further, except perhaps to magnify the paramecia still more and see inside. This, of course, is a subject for biology, but for the present we pass on and look still more closely at the water material itself, magnifying it two thousand times again. Now the drop of water extends about fifteen miles across, and if we look very closely at it we see a kind of teeming, something which no longer has a smooth appearance—it looks something like a crowd at a football game as seen from a very great distance. In order to see what this teeming is about, we will magnify it another two hundred and fifty times and we will see something similar to what is shown in Fig. 1-1. This is a picture of water magnified a billion times, but idealized in several ways.

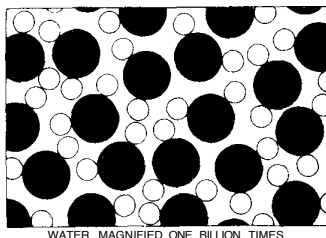


Figure 1-1

In the first place, the particles are drawn in a simple manner with sharp edges, which is inaccurate. Secondly, for simplicity, they are sketched almost schematically in a two-dimensional arrangement, but of course they are moving around in three dimensions. Notice that there are two kinds of "blobs" or circles to represent the atoms of oxygen (black) and hydrogen (white), and that each oxygen has two hydrogens tied to it. (Each little group of an oxygen with its two hydrogens is called a molecule.) The picture is idealized further in that the real particles in nature are continually jiggling and bouncing, turning and twisting around one another. You will have to imagine this as a dynamic rather than a static picture. Another thing that cannot be illustrated in a drawing is the fact that the particles are "stuck together"—that they attract each other, this one pulled by that one, etc. The whole group is "glued together," so to speak. On the other hand, the particles do not squeeze through each other. If you try to squeeze two of them too close together, they repel.

The atoms are 1 or 2×10^{-8} cm in radius. Now 10^{-8} cm is called an *angstrom* (just as another name), so we say they are 1 or 2 angstroms (\AA) in radius. Another way to remember their size is this: if an apple is magnified to the size of the earth, then the atoms in the apple are approximately the size of the original apple.

Now imagine this great drop of water with all of these jiggling particles stuck together and tagging along with each other. The water keeps its volume; it does not fall apart, because of the attraction of the molecules for each other. If the drop is on a slope, where it can move from one place to another, the water will flow, but it does not just disappear—things do not just fly apart—because of the molecular attraction. Now the jiggling motion is what we represent as *heat*: when we increase the temperature, we increase the motion. If we heat the water, the jiggling increases and the volume between the atoms increases, and if the heating continues there comes a time when the pull between the molecules is not enough to hold them together and they *do* fly apart and become separated from one another. Of course, this is how we manufacture steam out of water—by increasing the temperature; the particles fly apart because of the increased motion.

In Fig. 1-2 we have a picture of steam. This picture of steam fails in one respect: at ordinary atmospheric pressure there might be only a few molecules in a whole room, and there certainly would not be as many as three in this figure. Most squares this size would contain none—but we accidentally have two and a half or three in the picture (just so it would not be completely blank). Now in the case of steam we see the characteristic molecules more clearly than in the case of water. For simplicity, the molecules are drawn so that there is a 120° angle between them. In actual fact the angle is $105^\circ 3'$, and the distance between the center of a hydrogen and the center of the oxygen is 0.957 \AA , so we know this molecule very well.

Let us see what some of the properties of steam vapor or any other gas are. The molecules, being separated from one another, will bounce against the walls. Imagine a room with a number of tennis balls (a hundred or so) bouncing around in perpetual motion. When they bombard the wall, this pushes the wall away. (Of course we would have to push the wall back.) This means that the gas exerts a jittery force which our coarse senses (not being ourselves magnified a billion times) feels only as an *average push*. In order to confine a gas we must apply a pressure. Figure 1-3 shows a standard vessel for holding gases (used in all textbooks), a cylinder with a piston in it. Now, it makes no difference what the shapes of water molecules are, so for simplicity we shall draw them as tennis balls or little dots. These things are in perpetual motion in all directions. So many of them are hitting the top piston all the time that to keep it from being patiently knocked out of the tank by this continuous banging, we shall have to hold the piston down by a certain force, which we call the *pressure* (really, the pressure times the area is the force). Clearly, the force is proportional to the area, for if we increase the area but keep the number of molecules per cubic centimeter the same, we increase the number of collisions with the piston in the same proportion as the area was increased.

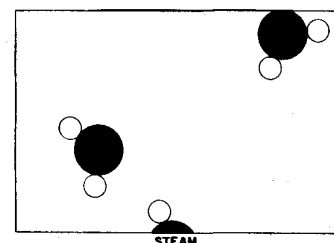


Figure 1-2

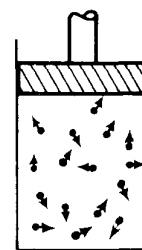


Figure 1-3

Now let us put twice as many molecules in this tank, so as to double the density, and let them have the same speed, i.e., the same temperature. Then, to a close approximation, the number of collisions will be doubled, and since each will be just as "energetic" as before, the pressure is proportional to the density. If we consider the true nature of the forces between the atoms, we would expect a slight decrease in pressure because of the attraction between the atoms, and a slight increase because of the finite volume they occupy. Nevertheless, to an excellent approximation, if the density is low enough that there are not many atoms, *the pressure is proportional to the density*.

We can also see something else: If we increase the temperature without changing the density of the gas, i.e., if we increase the speed of the atoms, what is going to happen to the pressure? Well, the atoms hit harder because they are moving faster, and in addition they hit more often, so the pressure increases. You see how simple the ideas of atomic theory are.

Let us consider another situation. Suppose that the piston moves inward, so that the atoms are slowly compressed into a smaller space. What happens when an atom hits the moving piston? Evidently it picks up speed from the collision. You can try it by bouncing a ping-pong ball from a forward-moving paddle, for example, and you will find that it comes off with more speed than that with which it struck. (Special example: if an atom happens to be standing still and the piston hits it, it will certainly move.) So the atoms are "hotter" when they come away from the piston than they were before they struck it. Therefore all the atoms which are in the vessel will have picked up speed. This means that *when we compress a gas slowly, the temperature of the gas increases*. So, under slow compression, a gas will increase in temperature, and under slow expansion it will decrease in temperature.

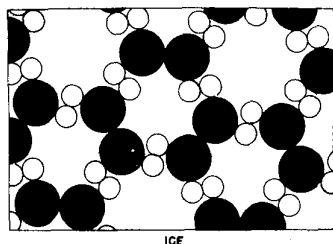


Figure 1-4

We now return to our drop of water and look in another direction. Suppose that we decrease the temperature of our drop of water. Suppose that the jiggling of the molecules of the atoms in the water is steadily decreasing. We know that there are forces of attraction between the atoms, so that after a while they will not be able to jiggle so well. What will happen at very low temperatures is indicated in Fig. 1-4: the molecules lock into a new pattern which is *ice*. This particular schematic diagram of ice is wrong because it is in two dimensions, but it is right qualitatively. The interesting point is that the material has a *definite place for every atom*, and you can easily appreciate that if somehow or other we were to hold all the atoms at one end of the drop in a certain arrangement, each atom in a certain place, then because of the structure of interconnections, which is rigid, the other end miles away (at our magnified scale) will have a definite location. So if we hold a needle of ice at one end, the other end resists our pushing it aside, unlike the case of water, in which the structure is broken down because of the increased jiggling so that the atoms all move around in different ways. The difference between solids and liquids is, then, that in a solid the atoms are arranged in some kind of an array, called a *crystalline array*, and they do not have a random position at long distances; the position of the atoms on one side of the crystal is determined by that of other atoms millions of atoms away on the other side of the crystal. Figure 1-4 is an invented arrangement for ice, and although it contains many of the correct features of ice, it is not the true arrangement. One of the correct features is that there is a part of the symmetry that is hexagonal. You can see that if we turn the picture around an axis by 120° , the picture returns to itself. So there is a *symmetry* in the ice which accounts for the six-sided appearance of snowflakes. Another thing we can see from Fig. 1-4 is why ice shrinks when it melts. The particular crystal pattern of ice shown here has many "holes" in it, as does the true ice structure. When the organization breaks down, these holes can be occupied by molecules. Most simple substances, with the exception of water and type metal, *expand* upon melting, because the atoms are closely packed in the solid crystal and upon melting need more room to jiggle around, but an open structure collapses, as in the case of water.

Now although ice has a "rigid" crystalline form, its temperature can change—ice has heat. If we wish, we can change the amount of heat. What is the heat in

the case of ice? The atoms are not standing still. They are jiggling and vibrating. So even though there is a definite order to the crystal—a definite structure—all of the atoms are vibrating "in place." As we increase the temperature, they vibrate with greater and greater amplitude, until they shake themselves out of place. We call this *melting*. As we decrease the temperature, the vibration decreases and decreases until, at absolute zero, there is a minimum amount of vibration that the atoms can have, but *not zero*. This minimum amount of motion that atoms can have is not enough to melt a substance, with one exception: helium. Helium merely decreases the atomic motions as much as it can, but even at absolute zero there is still enough motion to keep it from freezing. Helium, even at absolute zero, does not freeze, unless the pressure is made so great as to make the atoms squash together. If we increase the pressure, we *can* make it solidify.

1-3 Atomic processes

So much for the description of solids, liquids, and gases from the atomic point of view. However, the atomic hypothesis also describes *processes*, and so we shall now look at a number of processes from an atomic standpoint. The first process that we shall look at is associated with the surface of the water. What happens at the surface of the water? We shall now make the picture more complicated—and more realistic—by imagining that the surface is in air. Figure 1-5 shows the surface of water in air. We see the water molecules as before, forming a body of liquid water, but now we also see the surface of the water. Above the surface we find a number of things: First of all there are water molecules, as in steam. This is *water vapor*, which is always found above liquid water. (There is an equilibrium between the steam vapor and the water which will be described later.) In addition we find some other molecules—here two oxygen atoms stuck together by themselves, forming an *oxygen molecule*, there two nitrogen atoms also stuck together to make a nitrogen molecule. Air consists almost entirely of nitrogen, oxygen, some water vapor, and lesser amounts of carbon dioxide, argon, and other things. So above the water surface is the air, a gas, containing some water vapor. Now what is happening in this picture? The molecules in the water are always jiggling around. From time to time, one on the surface happens to be hit a little harder than usual, and gets knocked away. It is hard to see that happening in the picture because it is a *still* picture. But we can imagine that one molecule near the surface has just been hit and is flying out, or perhaps another one has been hit and is flying out. Thus, molecule by molecule, the water disappears—it evaporates. But if we *close* the vessel above, after a while we shall find a large number of molecules of water amongst the air molecules. From time to time, one of these vapor molecules comes flying down to the water and gets stuck again. So we see that what looks like a dead, uninteresting thing—a glass of water with a cover, that has been sitting there for perhaps twenty years—really contains a dynamic and interesting phenomenon which is going on all the time. To our eyes, our crude eyes, nothing is changing, but if we could see it a billion times magnified, we would see that from its own point of view it is always changing: molecules are leaving the surface, molecules are coming back.

Why do we see *no change*? Because just as many molecules are leaving as are coming back! In the long run "nothing happens." If we then take the top of the vessel off and blow the moist air away, replacing it with dry air, then the number of molecules leaving is just the same as it was before, because this depends on the jiggling of the water, but the number coming back is greatly reduced because there are so many fewer water molecules above the water. Therefore there are more going out than coming in, and the water evaporates. Hence, if you wish to evaporate water turn on the fan!

Here is something else: Which molecules leave? When a molecule leaves it is due to an accidental, extra accumulation of a little bit more than ordinary energy, which it needs if it is to break away from the attractions of its neighbors. Therefore, since those that leave have more energy than the average, the ones that are left have *less* average motion than they had before. So the liquid gradually

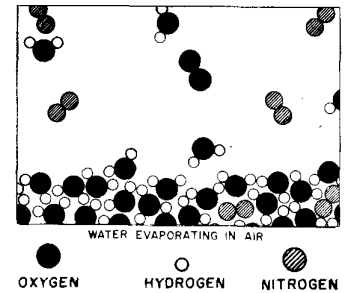


Figure 1-5

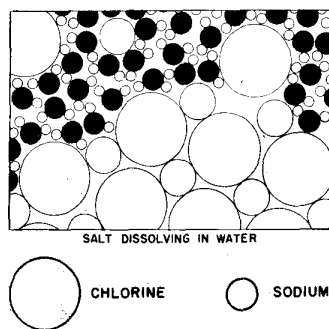


Figure 1-6

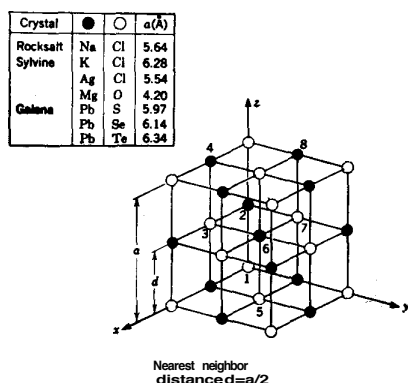


Figure 1-7

cools if it evaporates. Of course, when a molecule of vapor comes from the air to the water below there is a sudden great attraction as the molecule approaches the surface. This speeds up the incoming molecule and results in generation of heat. So when they leave they take away heat; when they come back they generate heat. Of course when there is no net evaporation the result is nothing—the water is not changing temperature. If we blow on the water so as to maintain a continuous preponderance in the number evaporating, then the water is cooled. Hence, blow on soup to cool it!

Of course you should realize that the processes just described are more complicated than we have indicated. Not only does the water go into the air, but also, from time to time, one of the oxygen or nitrogen molecules will come in and "get lost" in the mass of water molecules, and work its way into the water. Thus the air dissolves in the water; oxygen and nitrogen molecules will work their way into the water and the water will contain air. If we suddenly take the air away from the vessel, then the air molecules will leave more rapidly than they come in, and in doing so will make bubbles. This is very bad for divers, as you may know.

Now we go on to another process. In Fig. 1-6 we see, from an atomic point of view, a solid dissolving in water. If we put a crystal of salt in the water, what will happen? Salt is a solid, a crystal, an organized arrangement of "salt atoms." Figure 1-7 is an illustration of the three-dimensional structure of common salt, sodium chloride. Strictly speaking, the crystal is not made of atoms, but of what we call *ions*. An ion is an atom which either has a few extra electrons or has lost a few electrons. In a salt crystal we find chlorine ions (chlorine atoms with an extra electron) and sodium ions (sodium atoms with one electron missing). The ions all stick together by electrical attraction in the solid salt, but when we put them in the water we find, because of the attractions of the negative oxygen and positive hydrogen for the ions, that some of the ions jiggle loose. In Fig. 1-6 we see a chlorine ion getting loose, and other atoms floating in the water in the form of ions. This picture was made with some care. Notice, for example, that the hydrogen ends of the water molecules are more likely to be near the chlorine ion, while near the sodium ion we are more likely to find the oxygen end, because the sodium is positive and the oxygen end of the water is negative, and they attract electrically. Can we tell from this picture whether the salt is *dissolving in* water or *crystallizing out* of water? Of course we *cannot* tell, because while some of the atoms are leaving the crystal other atoms are rejoining it. The process is a *dynamic* one, just as in the case of evaporation, and it depends on whether there is more or less salt in the water than the amount needed for equilibrium. By equilibrium we mean that situation in which the rate at which atoms are leaving just matches the rate at which they are coming back. If there is almost no salt in the water, more atoms leave than return, and the salt dissolves. If, on the other hand, there are too many "salt atoms," more return than leave, and the salt is crystallizing.

In passing, we mention that the concept of a *molecule* of a substance is only approximate and exists only for a certain class of substances. It is clear in the case of water that the three atoms are actually stuck together. It is not so clear in the case of sodium chloride in the solid. There is just an arrangement of sodium and chlorine ions in a cubic pattern. There is no natural way to group them as "molecules of salt."

Returning to our discussion of solution and precipitation, if we increase the temperature of the salt solution, then the rate at which atoms are taken away is increased, and so is the rate at which atoms are brought back. It turns out to be very difficult, in general, to predict which way it is going to go, whether more or less of the solid will dissolve. Most substances dissolve more, but some substances dissolve less, as the temperature increases.

14 Chemical reactions

In all of the processes which have been described so far, the atoms and the ions have not changed partners, but of course there are circumstances in which the atoms do change combinations, forming new molecules. This is illustrated in 1-6

Fig. 1-8. A process in which the rearrangement of the atomic partners occurs is what we call a *chemical reaction*. The other processes so far described are called physical processes, but there is no sharp distinction between the two. (Nature does not care what we call it, she just keeps on doing it.) This figure is supposed to represent carbon burning in oxygen. In the case of oxygen, *two* oxygen atoms stick together very strongly. (Why do not *three* or *even four* stick together? That is one of the very peculiar characteristics of such atomic processes. Atoms are very special: they like certain particular partners, certain particular directions, and so on. It is the job of physics to analyze why each one wants what it wants. At any rate, two oxygen atoms form, saturated and happy, a molecule.)

The carbon atoms are supposed to be in a solid crystal (which could be graphite or diamond*). Now, for example, one of the oxygen molecules can come over to the carbon, and each atom can pick up a carbon atom and go flying off in a new combination—"carbon-oxygen"—which is a molecule of the gas called carbon monoxide. It is given the chemical name CO. It is very simple: the letters "CO" are practically a picture of that molecule. But carbon attracts oxygen much more than oxygen attracts oxygen or carbon attracts carbon. Therefore in this process the oxygen may arrive with only a little energy, but the oxygen and carbon will snap together with a tremendous vengeance and commotion, and everything near them will pick up the energy. A large amount of motion energy, kinetic energy, is thus generated. This of course is *burning*; we are getting *heat* from the combination of oxygen and carbon. The heat is ordinarily in the form of the molecular motion of the hot gas, but in certain circumstances it can be so enormous that it generates *light*. That is how one *gets flames*.

In addition, the carbon monoxide is not quite satisfied. It is possible for it to attach another oxygen, so that we might have a much more complicated reaction in which the oxygen is combining with the carbon, while at the same time there happens to be a collision with a carbon monoxide molecule. One oxygen atom could attach itself to the CO and ultimately form a molecule, composed of one carbon and two oxygens, which is designated CO₂ and called carbon dioxide. If we burn the carbon with very little oxygen in a very rapid reaction (for example, in an automobile engine, where the explosion is so fast that there is not time for it to make carbon dioxide) a considerable amount of carbon monoxide is formed. In many such rearrangements, a very large amount of energy is released, forming explosions, flames, etc., depending on the reactions. Chemists have studied these arrangements of the atoms, and found that every substance is some type of *arrangement of atoms*.

To illustrate this idea, let us consider another example. If we go into a field of small violets, we know what "that smell" is. It is some kind of *molecule*, or arrangement of atoms, that has worked its way into our noses. First of all, *how* did it work its way in? That is rather easy. If the smell is some kind of molecule in the air, jiggling around and being knocked every which way, it might have *accidentally* worked its way into the nose. Certainly it has no particular desire to get into our nose. It is merely one helpless part of a jostling crowd of molecules, and in its aimless wanderings this particular chunk of matter happens to find itself in the nose.

Now chemists can take special molecules like the odor of violets, and analyze them and tell us the *exact arrangement* of the atoms in space. We know that the carbon dioxide molecule is straight and symmetrical: O—C—O. (That can be determined easily, too, by physical methods.) However, even for the vastly more complicated arrangements of atoms that there are in chemistry, one can, by a long, remarkable process of detective work, find the arrangements of the atoms. Figure 1-9 is a picture of the air in the neighborhood of a violet; again we find nitrogen and oxygen in the air, and water vapor. (Why is there water vapor? Because the violet is *wet*. All plants transpire.) However, we also see a "monster" composed of carbon atoms, hydrogen atoms, and oxygen atoms, which have picked a certain particular pattern in which to be arranged. It is a much more complicated arrange-

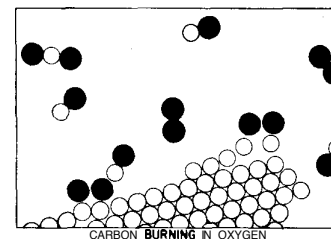


Figure 1-8

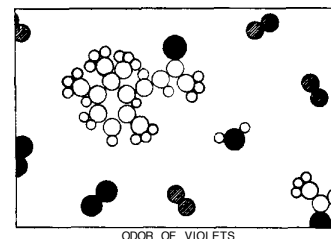


Figure 1-9

*One can burn a diamond in air.

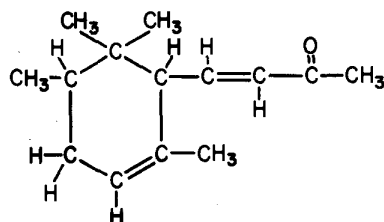


Fig. 1-10. The substance pictured is α -ionone.

ment than that of carbon dioxide; in fact, it is an enormously complicated arrangement. Unfortunately, we cannot picture all that is really known about it chemically, because the precise arrangement of all the atoms is actually known in three dimensions, while our picture is in only two dimensions. The six carbons which form a ring do not form a flat ring, but a kind of "puckered" ring. All of the angles and distances are known. So a chemical *formula* is merely a picture of such a molecule. When the chemist writes such a thing on the blackboard, he is trying to "draw," roughly speaking, in two dimensions. For example, we see a "ring" of six carbons, and a "chain" of carbons hanging on the end, with an oxygen second from the end, three hydrogens tied to that carbon, two carbons and three hydrogens sticking up here, etc.

How does the chemist find what the arrangement is? He mixes bottles full of stuff together, and if it turns red, it tells him that it consists of one hydrogen and two carbons tied on here; if it turns blue, on the other hand, that is not the way it is at all. This is one of the most fantastic pieces of detective work that has ever been done—organic chemistry. To discover the arrangement of the atoms in these enormously complicated arrays the chemist looks at what happens when he mixes two different substances together. The physicist could never quite believe that the chemist knew what he was talking about when he described the arrangement of the atoms. For about twenty years it has been possible, in some cases, to look at such molecules (not quite as complicated as this one, but some which contain parts of it) by a physical method, and it has been possible to locate every atom, not by looking at colors, but by *measuring where they are*. And lo and behold!, the chemists are almost always correct.

It turns out, in fact, that in the odor of violets there are three slightly different molecules, which differ only in the arrangement of the hydrogen atoms.

One problem of chemistry is to name a substance, so that we will know what it is. Find a name for this shape! Not only must the name tell the shape, but it must also tell that there is an oxygen atom, there a hydrogen—exactly what and where each atom is. So we can appreciate that the chemical names must be complex in order to be complete. You see that the name of this thing in the more complete form that will tell you the structure of it is 4-(2, 2, 3, 6 tetramethyl-5-cyclohexenyl)-3-buten-2-one, and that tells you that this is the arrangement. We can appreciate the difficulties that the chemists have, and also appreciate the reason for such long names. It is not that they wish to be obscure, but they have an extremely difficult problem in trying to describe the molecules in words!

How do we *know* that there are atoms? By one of the tricks mentioned earlier: we make the *hypothesis* that there are atoms, and one after the other results come out the way we predict, as they ought to if things *are* made of atoms. There is also somewhat more direct evidence, a good example of which is the following: The atoms are so small that you cannot see them with a light microscope—in fact, not even with an *electron* microscope. (With a light microscope you can only see things which are much bigger.) Now if the atoms are always in motion, say in water, and we put a big ball of something in the water, a ball much bigger than the atoms, the ball will jiggle around—much as in a push ball game, where a great big ball is pushed around by a lot of people. The people are pushing in various directions, and the ball moves around the field in an irregular fashion. So, in the same way, the "large ball" will move because of the inequalities of the collisions on one side to the other, from one moment to the next. Therefore, if we look at very tiny particles (colloids) in water through an excellent microscope, we see a perpetual jiggling of the particles, which is the result of the bombardment of the atoms. This is called the *Brownian motion*.

We can see further evidence for atoms in the structure of crystals. In many cases the structures deduced by x-ray analysis agree in their spatial "shapes" with the forms actually exhibited by crystals as they occur in nature. The angles between the various "faces" of a crystal agree, within seconds of arc, with angles deduced on the assumption that a crystal is made of many "layers" of atoms.

Everything is made of atoms. That is the key hypothesis. The most important hypothesis in all of biology, for example, is that *everything that animals do, atoms*

do. In other words, *there is nothing that living things do that cannot be understood from the point of view that they are made of atoms acting according to the laws of physics*. This was not known from the beginning: it took some experimenting and theorizing to suggest this hypothesis, but now it is accepted, and it is the most useful theory for producing new ideas in the field of biology.

If a piece of steel or a piece of salt, consisting of atoms one next to the other, can have such interesting properties; if water—which is nothing but these little blobs, mile upon mile of the same thing over the earth—can form waves and foam, and make rushing noises and strange patterns as it runs over cement; if all of this, all the life of a stream of water, can be nothing but a pile of atoms, *how much more is possible?* If instead of arranging the atoms in some definite pattern, again and again repeated, on and on, or even forming little lumps of complexity like the odor of violets, we make an arrangement which is *always different* from place to place, with different kinds of atoms arranged in many ways, continually changing, not repeating, how much more marvelously is it possible that this thing might behave? Is it possible that that "thing" walking back and forth in front of you, talking to you, is a great glob of these atoms in a very complex arrangement, such that the sheer complexity of it staggers the imagination as to what it can do? When we say we are a pile of atoms, we do not mean we are *merely* a pile of atoms, because a pile of atoms which is not repeated from one to the other might well have the possibilities which you see before you in the mirror.

Basic Physics

2-1 Introduction

In this chapter, we shall examine the most fundamental ideas that we have about physics—the nature of things as we see them at the present time. We shall not discuss the history of how we know that all these ideas are true; you will learn these details in due time.

The things with which we concern ourselves in science appear in myriad forms, and with a multitude of attributes. For example, if we stand on the shore and look at the sea, we see the water, the waves breaking, the foam, the sloshing motion of the water, the sound, the air, the winds and the clouds, the sun and the blue sky, and light; there is sand and there are rocks of various hardness and permanence, color and texture. There are animals and seaweed, hunger and disease, and the observer on the beach; there may be even happiness and thought. Any other spot in nature has a similar variety of things and influences. It is always as complicated as that, no matter where it is. Curiosity demands that we ask questions, that we try to put things together and try to understand this multitude of aspects as perhaps resulting from the action of a relatively small number of elemental things and forces acting in an infinite variety of combinations.

For example: Is the sand other than the rocks? That is, is the sand perhaps nothing but a great number of very tiny stones? Is the moon a great rock? If we understood rocks, would we also understand the sand and the moon? Is the wind a sloshing of the air analogous to the sloshing motion of the water in the sea? What common features do different movements have? What is common to different kinds of sound? How many different colors are there? And so on. In this way we try gradually to analyze all things, to put together things which at first sight look different, with the hope that we may be able to *reduce* the number of *different* things and thereby understand them better.

A few hundred years ago, a method was devised to find partial answers to such questions. *Observation*, *reason*, and *experiment* make up what we call the *scientific method*. We shall have to limit ourselves to a bare description of **our** basic view of what is sometimes called *fundamental physics*, or fundamental ideas which have arisen from the application of the scientific method.

What do we mean by "understanding" something? We can imagine that this complicated array of moving things which constitutes "the world" is something like a great chess game being played by the gods, and we are observers of the game. We do not know what the rules of the game are; all we are allowed to do is to *watch* the playing. Of course, if we watch long enough, we may eventually catch on to a few of the rules. *The rules of the game* are what we mean by *fundamental physics*. Even if we knew every rule, however, we might not be able to understand why a particular move is made in the game, merely because it is too complicated and our minds are limited. If you play chess you must know that it is easy to learn all the rules, and yet it is often very hard to select the best move or to understand why a player moves as he does. So it is in nature, only much more so; but we may be able at least to find all the rules. Actually, we do not have all the rules now. (Every once in a while something like castling is going on that we still do not understand.) Aside from not knowing all of the rules, what we really can explain in terms of those rules is very limited, because almost all situations are so enormously complicated that we cannot follow the plays of the game using the rules, much less tell what is going to happen next. We must, therefore, limit ourselves to the more basic question of the rules of the game. If we know the rules, we consider that we "understand" the world.

2-1 Introduction

2-2 Physics before 1920

2-3 Quantum physics

2-4 Nuclei and particles

How can we tell whether the rules which we "guess" at are really right if we cannot analyze the game very well? There are, roughly speaking, three ways. First, there may be situations where nature has arranged, or we arrange nature, to be simple and to have so few parts that we can predict exactly what will happen, and thus we can check how our rules work. (In one corner of the board there may be only a few chess pieces at work, and that we can figure out exactly.)

A second good way to check rules is in terms of less specific rules derived from them. For example, the rule on the move of a bishop on a chessboard is that it moves only on the diagonal. One can deduce, no matter how many moves may be made, that a certain bishop will always be on a red square. So, without being able to follow the details, we can always check our idea about the bishop's motion by finding out whether it is always on a red square. Of course it will be, for a long time, until all of a sudden we find that it is on a *black* square (what happened of course, is that in the meantime it was captured, another pawn crossed for queening, and it turned into a bishop on a black square). That is the way it is in physics. For a long time we will have a rule that works excellently in an over-all way, even when we cannot follow the details, and then some time we may discover a *new rule*. From the point of view of basic physics, the most interesting phenomena are of course in the *new* places, the places where the rules do not work—not the places where they *do* work! That is the way in which we discover new rules.

The third way to tell whether our ideas are right is relatively crude but probably the most powerful of them all. That is, by rough *approximation*. While we may not be able to tell why Alekhine moves *this particular piece*, perhaps we can *roughly* understand that he is gathering his pieces around the king to protect it, more or less, since that is the sensible thing to do in the circumstances. In the same way, we can often understand nature, more or less, without being able to see what *every little piece* is doing, in terms of our understanding of the game.

At first the phenomena of nature were roughly divided into classes, like heat, electricity, mechanics, magnetism, properties of substances, chemical phenomena, light or optics, x-rays, nuclear physics, gravitation, meson phenomena, etc. However, the aim is to see *complete nature* as different aspects of *one set* of phenomena. That is the problem in basic theoretical physics, today—to *find the laws behind experiment*; to *amalgamate these classes*. Historically, we have always been able to amalgamate them, but as time goes on new things are found. We were amalgamating very well, when all of a sudden x-rays were found. Then we amalgamated some more, and mesons were found. Therefore, at any stage of the game, it always looks rather messy. A great deal is amalgamated, but there are always many wires or threads hanging out in all directions. That is the situation today, which we shall try to describe.

Some historic examples of amalgamation are the following. First, take *heat* and *mechanics*. When atoms are in motion, the more motion, the more heat the system contains, and so *heat and all temperature effects can be represented by the laws of mechanics*. Another tremendous amalgamation was the discovery of the relation between electricity, magnetism, and light, which were found to be different aspects of the same thing, which we call today the *electromagnetic field*. Another amalgamation is the unification of chemical phenomena, the various properties of various substances, and the behavior of atomic particles, which is in the *quantum mechanics of chemistry*.

The question is, of course, is it going to be possible to amalgamate *everything*, and merely discover that this world represents different aspects of *one* thing? Nobody knows. All we know is that as we go along, we find that we can amalgamate pieces, and then we find some pieces that do not fit, and we keep trying to put the jigsaw puzzle together. Whether there are a finite number of pieces, and whether there is even a border to the puzzle, is of course unknown. It will never be known until we finish the picture, if ever. What we wish to do here is to see to what extent this amalgamation process has gone on, and what the situation is at present, in understanding basic phenomena in terms of the smallest set of principles. To express it in a simple manner, *what are things made of and how few elements are there?*

2-2 Physics before 1920

It is a little difficult to begin at once with the present view, so we shall first see how things looked in about 1920 and then take a few things out of that picture. Before 1920, our world picture was something like this: The "stage" on which the universe goes is the three-dimensional *space* of geometry, as described by Euclid, and things change in a medium called *time*. The elements on the stage are *particles*, for example the atoms, which have some *properties*. First, the property of inertia: if a particle is moving it keeps on going in the same direction unless *forces* act upon it. The second element, then, is *forces*, which were then thought to be of two varieties: First, an enormously complicated, detailed kind of interaction force which held the various atoms in different combinations in a complicated way, which determined whether salt would dissolve faster or slower when we raise the temperature. The other force that was known was a long-range interaction—a smooth and quiet attraction—which varied inversely as the square of the distance, and was called *gravitation*. This law was known and was very simple. *Why* things remain in motion when they are moving, or *why* there is a law of gravitation was, of course, not known.

A description of nature is what we are concerned with here. From this point of view, then, a gas, and indeed *all* matter, is a myriad of moving particles. Thus many of the things we saw while standing at the seashore can immediately be connected. First the pressure: this comes from the collisions of the atoms with the walls or whatever; the drift of the atoms, if they are all moving in one direction on the average, is wind; the *random* internal motions are the *heat*. There are waves of excess density, where too many particles have collected, and so as they tush off they push up piles of particles farther out, and so on. This wave of excess density is *sound*. It is a tremendous achievement to be able to understand so much. Some of these things were described in the previous chapter.

What *kinds* of particles are there? There were considered to be 92 at that time: 92 different kinds of atoms were ultimately discovered. They had different names associated with their chemical properties.

The next part of the problem was, *what are the short-range forces?* Why does carbon attract one oxygen or perhaps two oxygens, but not three oxygens? What is the machinery of interaction between atoms? Is it gravitation? The answer is no. Gravity is entirely too weak. But imagine a force analogous to gravity, varying inversely with the square of the distance, but enormously more powerful and having one difference. In gravity everything attracts everything else, but now imagine that there are *two kinds* of "things," and that this new force (which is the electrical force, of course) has the property that likes *repel* but unlikes *attract*. The "thing" that carries this strong interaction is called *charge*.

Then what do we have? Suppose that we have two unlikes that attract each other, a plus and a minus, and that they stick very close together. Suppose we have another charge some distance away. Would it feel any attraction? It would feel *practically none*, because if the first two are equal in size, the attraction for the one and the repulsion for the other balance out. Therefore there is very little force at any appreciable distance. On the other hand, if we get *very close* with the extra charge, *attraction* arises, because the repulsion of likes and attraction of unlikes will tend to bring unlikes closer together and push likes farther apart. Then the repulsion will be *less* than the attraction. This is the reason why the atoms, which are constituted out of plus and minus electric charges, feel very little force when they are separated by appreciable distance (aside from gravity). When they come close together, they can "see inside" each other and rearrange their charges, with the result that they have a very strong interaction. The ultimate basis of an interaction between the atoms is *electrical*. Since this force is so enormous, all the plusses and all minuses will normally come together in as intimate a combination as they can. All things, even ourselves, are made of fine-grained, enormously strongly interacting plus and minus parts, all neatly balanced out. Once in a while, by accident, we may rub off a few minuses or a few plusses (usually it is easier to rub off minuses), and in those circumstances we find the force of electricity *unbalanced*, and we can then see the effects of these electrical attractions.

To give an idea of how much stronger electricity is than gravitation, consider two grains of sand, a millimeter across, thirty meters apart. If the force between them were not balanced, if everything attracted everything else instead of likes repelling, so that there were no cancellation, how much force would there be? There would be a force of *three million tons* between the two! You see, there is very, very little excess or deficit of the number of negative or positive charges necessary to produce appreciable electrical effects. This is, of course, the reason why you cannot see the difference between an electrically charged or uncharged thing—so few particles are involved that they hardly make a difference in the weight or size of an object.

With this picture the atoms were easier to understand. They were thought to have a "nucleus" at the center, which is positively electrically charged and very massive, and the nucleus is surrounded by a certain number of "electrons" which are very light and negatively charged. Now we go a little ahead in our story to remark that in the nucleus itself there were found two kinds of particles, protons and neutrons, almost of the same weight and very heavy. The protons are electrically charged and the neutrons are neutral. If we have an atom with six protons inside its nucleus, and this is surrounded by six electrons (the negative particles in the ordinary world of matter are all electrons, and these are very light compared with the protons and neutrons which make nuclei), this would be atom number six in the chemical table, and it is called carbon. Atom number eight is called oxygen, etc., because the chemical properties depend upon the electrons on the *outside*, and in fact only upon *how many* electrons there are. So the *chemical* properties of a substance depend only on a number, the number of electrons. (The whole list of elements of the chemists really could have been called 1, 2, 3, 4, 5, etc. Instead of saying "carbon," we could say "element six," meaning six electrons, but of course, when the elements were first discovered, it was not known that they could be numbered that way, and secondly, it would make everything look rather complicated. It is better to have names and symbols for these things, rather than to call everything by number.)

More was discovered about the electrical force. The natural interpretation of electrical interaction is that two objects simply attract each other: plus against minus. However, this was discovered to be an inadequate idea to represent it. A more adequate representation of the situation is to say that the existence of the positive charge, in some sense, distorts, or creates a "condition" in space, so that when we put the negative charge in, it feels a force. This potentiality for producing a force is called an *electric field*. When we put an electron in an electric field, we say it is "pulled." We then have two rules: (a) charges make a field, and (b) charges in fields have forces on them and move. The reason for this will become clear when we discuss the following phenomena: If we were to charge a body, say a comb, electrically, and then place a charged piece of paper at a distance and move the comb back and forth, the paper will respond by always pointing to the comb. If we shake it faster, it will be discovered that the paper is a little behind, *there is a delay* in the action. (At the first stage, when we move the comb rather slowly, we find a complication which is *magnetism*. Magnetic influences have to do with *charges in relative motion*, so magnetic forces and electric forces can really be attributed to one field, as two different aspects of exactly the same thing. A changing electric field cannot exist without magnetism.) If we move the charged paper farther out, the delay is greater. Then an interesting thing is observed. Although the forces between two charged objects should go inversely as the *square* of the distance, it is found, when we shake a charge, that the influence extends *very much farther out* than we would guess at first sight. That is, the effect falls off more slowly than the inverse square.

Here is an analogy: If we are in a pool of water and there is a floating cork very close by, we can move it "directly" by pushing the water with another cork. If you looked only at the two *corks*, all you would see would be that one moved immediately in response to the motion of the other—there is some kind of "*inter-action*" between them. Of course, what we really do is to disturb the *water*; the *water* then disturbs the other cork. We could make up a "law" that if you pushed

the water a little bit, an object close by in the water would move. If it were farther away, of course, the second cork would scarcely move, for we move the water *locally*. On the other hand, if we jiggle the cork a new phenomenon is involved, in which the motion of the water moves the water there, etc., and *waves* travel away, so that by jiggling, there is an influence *wry much farther out*, an oscillatory influence, that cannot be understood from the direct interaction. Therefore the idea of direct interaction must be replaced with the existence of the water, or—in the electrical case, with what we call the *electromagnetic field*.

The electromagnetic field can carry waves; some of these waves are *light*, others are used in *radio broadcasts*, but the general name is *electromagnetic waves*. These oscillatory waves can have various *frequencies*. The only thing that is really different from one wave to another is the *frequency of oscillation*. If we shake a charge back and forth more and more rapidly, and look at the effects, we get a whole series of different kinds of effects, which are all unified by specifying but one number, the number of oscillations per second. The usual "pickup" that we get from electric currents in the circuits in the walls of a building have a frequency of about one hundred cycles per second. If we increase the frequency to 500 or 1000 kilocycles (1 kilocycle = 1000 cycles) per second, we are "on the air," for this is the frequency range which is used for radio broadcasts. (Of course it has nothing to do with the *air*! We can have radio broadcasts without any air.) If we again increase the frequency, we come into the range that is used for FM and TV. Going still further, we use certain short waves, for example for *radar*. Still higher, and we do not need an instrument to "see" the stuff, we can see it with the human eye. In the range of frequency from 5×10^{14} to 5×10^{15} cycles per second our eyes would see the oscillation of the charged comb, if we could shake it that fast, as red, blue, or violet light, depending on the frequency. Frequencies below this range are called infrared, and above it, ultraviolet. The fact that we can see in a particular frequency range makes that part of the electromagnetic spectrum no more impressive than the other parts from a physicist's standpoint, but from a human standpoint, of course, it *is* more interesting. If we go up even higher in frequency, we get x-rays. X-rays are nothing but very high-frequency light. If we go still higher, we get gamma rays. These two terms, x-rays and gamma rays, are used almost synonymously. Usually electromagnetic rays coming from nuclei are called gamma rays, while those of high energy from atoms are called x-rays, but at the same frequency they are indistinguishable physically, no matter what their source. If we go to still higher frequencies, say to 10^{24} cycles per second, we find that we can make those waves artificially, for example with the synchrotron here at Caltech. We can find electromagnetic waves with stupendously high frequencies—with even a thousand times more rapid oscillation—in the waves found in *cosmic rays*. These waves cannot be controlled by us.

Table 2-1
The Electromagnetic Spectrum

Frequency in oscillations/sec	Name	Rough behavior
10^2	Electrical disturbance	Field
$5 \times 10^5 - 10^6$	Radio broadcast	} Waves
10^8	FM — TV	
10^{10}	Radar	
$5 \times 10^{14} - 10^{15}$	Light	
10^{18}	X-rays	} Particle
10^{21}	γ -rays, nuclear	
10^{24}	γ -rays, "artificial"	
10^{27}	γ -rays, in cosmic rays	

2-3 Quantum physics

Having described the idea of the electromagnetic field, and that this field can carry waves, we soon learn that these waves actually behave in a strange way which seems very unwavelike. At higher frequencies they behave much more like *particles*! It is *quantum mechanics*, discovered just after 1920, which explains this strange behavior. In the years before 1920, the picture of space as a three-dimensional space, and of time as a separate thing, was changed by Einstein, first into a combination which we call space-time, and then still further into a *curved* space-time to represent gravitation. So the "stage" is changed into space-time, and gravitation is presumably a modification of space-time. Then it was also found that the rules for the motions of particles were incorrect. The mechanical rules of "inertia" and "forces" are *wrong*—Newton's laws are *wrong*—in the world of atoms. Instead, it was discovered that things on a small scale behave *nothing like* things on a large scale. That is what makes physics difficult—and very interesting. It is hard because the way things behave on a small scale is so "unnatural"; we have no direct experience with it. Here things behave like nothing we know of, so that it is impossible to describe this behavior in any other than analytic ways. It is difficult, and takes a lot of imagination.

Quantum mechanics has many aspects. In the first place, the idea that a particle has a definite location and a definite speed is no longer allowed; that is wrong. To give an example of how wrong classical physics is, there is a rule in quantum mechanics that says that one cannot know both where something is and how fast it is moving. The uncertainty of the momentum and the uncertainty of the position are complementary, and the product of the two is constant. We can write the law like this: $\Delta x \Delta p \geq \hbar/2$, but we shall explain it in more detail later. This rule is the explanation of a very mysterious paradox: if the atoms are made out of plus and minus charges, why don't the minus charges simply sit on top of the plus charges (they attract each other) and get so close as to completely cancel them out? *Why are atoms so big?* Why is the nucleus at the center with the electrons around it? It was first thought that this was because the nucleus was so big; but no, the nucleus is *very small*. An atom has a diameter of about 10^{-8} cm. The nucleus has a diameter of about 10^{-13} cm. If we had an atom and wished to see the nucleus, we would have to magnify it until the whole atom was the size of a large room, and then the nucleus would be a bare speck which you could just about make out with the eye, but very nearly *all the weight* of the atom is in that infinitesimal *nucleus*. What keeps the electrons from simply falling in? This principle: If they were in the nucleus, we would know their position precisely, and the uncertainty principle would then require that they have a very *large* (but uncertain) momentum, i.e., a very large *kinetic energy*. With this energy they would break away from the nucleus. They make a compromise: they leave themselves a little room for this uncertainty and then jiggle with a certain amount of minimum motion in accordance with this rule. (Remember that when a crystal is cooled to absolute zero, we said that the atoms do not stop moving, they still jiggle. Why? If they stopped moving, we would know where they were and that they had zero motion, and that is against the uncertainty principle. We cannot know where they are and how fast they are moving, so they must be continually wiggling in there!)

Another most interesting change in the ideas and philosophy of science brought about by quantum mechanics is this: it is not possible to predict *exactly* what will happen in any circumstance. For example, it is possible to arrange an atom which is ready to emit light, and we can measure when it has emitted light by picking up a photon particle, which we shall describe shortly. We cannot, however, predict *when* it is going to emit the light or, with several atoms, *which one* is going to. You may say that this is because there are some internal "wheels" which we have not looked at closely enough. No, there *are* no internal wheels; nature, as we understand it today, behaves in such a way that *it is fundamentally impossible* to make a precise prediction of *exactly what will happen* in a given experiment. This is a horrible thing; in fact, philosophers have said before that one of the fundamental requisites of science is that whenever you set **up** the same

conditions, the same thing must happen. This is simply *not true*, it is *not* a fundamental condition of science. The fact is that the same thing does not happen, that we can find only an average, statistically, **as to** what happens. Nevertheless, science has not completely collapsed. Philosophers, incidentally, say a great deal about what is *absolutely necessary* for science, and it is always, so far as one can see, rather naive, and probably wrong. For example, some philosopher or other said it is fundamental to the scientific effort that if an experiment is performed in, say, Stockholm, and then the same experiment is done in, say, Quito, the *same results* must occur. That is quite false. It is not necessary that *science* do that; it may be a *fact of experience*, but it is not necessary. For example, if one of the experiments is to look out at the sky and see the aurora borealis in Stockholm, you do not see it in Quito; that is a different phenomenon. "But," you say, "that is something that has to do with the outside; can you close yourself up in a box in Stockholm and pull down the shade and get any difference?" Surely. If we take a pendulum on a universal joint, and pull it out and let go, then the pendulum will swing almost in a plane, but not quite. Slowly the plane keeps changing in Stockholm, but not in Quito. The blinds are down, too. The fact that this happened does not bring on the destruction of science. What *is* the fundamental hypothesis of science, the fundamental philosophy? We stated it in the first chapter: *the sole test of the validity of any idea is experiment*. If it turns out that most experiments work out the same in Quito as they do in Stockholm, then those "most experiments" will be used to formulate some general law, and those experiments which do not come out the same we will say were a result of the environment near Stockholm. We will invent some way to summarize the results of the experiment, and we do not have to be told ahead of time what this way will look like. If we are told that the same experiment will always produce the same result, that is all very well, but if when we try it, it does *not*, then it does *not*. We just have to take what we see, and then formulate all the rest of our ideas in terms of our actual experience.

Returning again to quantum mechanics and fundamental physics, we cannot go into details of the quantum-mechanical principles at this time, of course, because these are rather difficult to understand. We shall assume that they are there, and go on to describe what some of the consequences are. One of the consequences is that things which we used to consider as waves also behave like particles, and particles behave like waves; in fact everything behaves the same way. There is no distinction between a wave and a particle. So quantum mechanics *unifies* the idea of the field and its waves, and the particles, all into one. Now it is true that when the frequency is low, the field aspect of the phenomenon is more evident, or more useful as an approximate description in terms of everyday experiences. But as the frequency increases, the particle aspects of the phenomenon become more evident with the equipment with which we usually make the measurements. In fact, although we mentioned many frequencies, no phenomenon directly involving a frequency has yet been detected above approximately 10^{12} cycles per second. We only *deduce* the higher frequencies from the energy of the particles, by a rule which assumes that the particle-wave idea of quantum mechanics is valid.

Thus we have a new view of electromagnetic interaction. We have a new kind of *particle* to add to the electron, the proton, and the neutron. That new particle is called a *photon*. The new view of the interaction of electrons and protons that is electromagnetic theory, but with everything quantum-mechanically correct, is called *quantum electrodynamics*. This fundamental theory of the interaction of light and matter, or electric field and charges, is our greatest success so far in physics. In this one theory we have the basic rules for all ordinary phenomena except for gravitation and nuclear processes. For example, out of quantum electrodynamics come all known electrical, mechanical, and chemical laws: the laws for the collision of billiard balls, the motions of wires in magnetic fields, the specific heat of carbon monoxide, the color of neon signs, the density of salt, and the reactions of hydrogen and oxygen to make water are all consequences of this one law. All these details can be worked out if the situation is simple enough for us to make an approximation, which is almost never, but often we can understand more

or less what is happening. At the present time no exceptions are found to the quantum-electrodynamic laws outside the nucleus, and there we do not know whether there is an exception because we simply do not know what is going on in the nucleus.

In principle, then, quantum electrodynamics is the theory of all chemistry, and of life, if life is ultimately reduced to chemistry and therefore just to physics because chemistry is already reduced (the part of physics which is involved in chemistry being already known). Furthermore, the same quantum electrodynamics, this great thing, predicts a lot of new things. In the first place, it tells the properties of very high-energy photons, gamma rays, etc. It predicted another very remarkable thing: besides the electron, there should be another particle of the same mass, but of opposite charge, called a *positron*, and these two, coming together, could annihilate each other with the emission of light or gamma rays. (After all, light and gamma rays are all the same, they are just different points on a frequency scale.) The generalization of this, that for each particle there is an antiparticle, turns out to be true. In the case of electrons, the antiparticle has another name—it is called a positron, but for most other particles, it is called anti-so-and-so, like antiproton or antineutron. In quantum electrodynamics, *two numbers* are put in and most of the other numbers in the world are supposed to come out. The two numbers that are put in are called the mass of the electron and the charge of the electron. Actually, that is not quite true, for we have a whole set of numbers for chemistry which tells how heavy the nuclei are. That leads us to the next part.

2-4 Nuclei and particles

What are the nuclei made of, and how are they held together? It is found that the nuclei are held together by enormous forces. When these are released, the energy released is tremendous compared with chemical energy, in the same ratio as the atomic bomb explosion is to a TNT explosion, because, of course, the atomic bomb has to do with changes inside the nucleus, while the explosion of TNT has to do with the changes of the electrons on the outside of the atoms. The question is, what are the forces which hold the protons and neutrons together in the nucleus? Just as the electrical interaction can be connected to a particle, a photon, Yukawa suggested that the forces between neutrons and protons also have a field of some kind, and that when this field jiggles it behaves like a particle. Thus there could be some other particles in the world besides protons and neutrons, and he was able to deduce the properties of these particles from the already known characteristics of nuclear forces. For example, he predicted they should have a mass of two or three hundred times that of an electron; and lo and behold, in cosmic rays there was discovered a particle of the right mass! But it later turned out to be the wrong particle. It was called a m-meson, or muon.

However, a little while later, in 1947 or 1948, another particle was found, the p-meson, or pion, which satisfied Yukawa's criterion. Besides the proton and the neutron, then, in order to get nuclear forces we must add the pion. Now, you say, "Oh great!, with this theory we make quantum nucleodynamics using the pions just like Yukawa wanted to do, and see if it works, and everything will be explained." Bad luck. It turns out that the calculations that are involved in this theory are so difficult that no one has ever been able to figure out what the consequences of the theory are, or to check it against experiment, and this has been going on now for almost twenty years!

So we are stuck with a theory, and we do not know whether it is right or wrong, but we do know that it is a *little* wrong, or at least incomplete. While we have been dawdling around theoretically, trying to calculate the consequences of this theory, the experimentalists have been discovering some things. For example, they had already discovered this m-meson or muon, and we do not yet know where it fits. Also, in cosmic rays, a large number of other "extra" particles were found. It turns out that today we have approximately thirty particles, and it is very difficult to understand the relationships of all these particles, and what nature

, wants them for, or what the connections are from one to another. We do not today understand these various particles as different aspects of the same thing, and the fact that we have so many unconnected particles is a representation of the fact that we have so much unconnected information without a good theory. After the great successes of quantum electrodynamics, there is a certain amount of knowledge of nuclear physics which is rough knowledge, sort of half experience and half theory, assuming a type of force between protons and neutrons and seeing what will happen, but not really understanding where the force comes from. Aside from that, we have made very little progress. We have collected an enormous number of chemical elements. In the chemical case, there suddenly appeared a relationship among these elements which was unexpected, and which is embodied in the periodic table of Mendeleev. For example, sodium and potassium are about the same in their chemical properties and are found in the same column in the Mendeleev chart. We have been seeking a Mendeleev-type chart for the new particles. One such chart of the new particles was made independently by Gell-Mann in the U.S.A. and Nishijima in Japan. The basis of their classification is a new number, like the electric charge, which can be assigned to each particle, called its "strangeness," S . This number is conserved, like the electric charge, in reactions which take place by nuclear forces.

In Table 2-2 are listed all the particles. We cannot discuss them much at this stage, but the table will at least show you how much we do not know. Underneath each particle its mass is given in a certain unit, called the Mev. One Mev is equal to 1.782×10^{-27} gram. The reason this unit was chosen is historical, and we shall not go into it now. More massive particles are put higher up on the chart; we see that a neutron and a proton have almost the same mass. In vertical columns we have put the particles with the same electrical charge, all neutral objects in one column, all positively charged ones to the right of this one, and all negatively charged objects to the left.

Particles are shown with a solid line and "resonances" with a dashed one. Several particles have been omitted from the table. These include the important zero-mass, zero-charge particles, the photon and the graviton, which do not fall into the baryon-meson-lepton classification scheme, and also some of the newer resonances (K^* , ρ , π). The antiparticles of the mesons are listed in the table, but the antiparticles of the leptons and baryons would have to be listed in another table which would look exactly like this one reflected on the zero-charge column. Although all of the particles except the electron, neutrino, photon, graviton, and proton are unstable, decay products have been shown only for the resonances. Strangeness assignments are not applicable for leptons, since they do not interact strongly with nuclei.

All particles which are together with the neutrons and protons are called *baryons*, and the following ones exist: There is a "lambda," with a mass of 1154 Mev, and three others, called sigmas, minus, neutral, and plus, with several masses almost the same. There are groups or multiplets with almost the same mass, within one or two percent. Each particle in a multiplet has the same strangeness. The first multiplet is the proton-neutron doublet, and then there is a singlet (the lambda) then the sigma triplet, and finally the xi doublet. Very recently, in 1961, even a few more particles were found. Or *are* they particles? They live so short a time, they disintegrate almost instantaneously, as soon as they are formed, that we do not know whether they should be considered as new particles, or some kind of "resonance" interaction of a certain definite energy between the A and T products into which they disintegrate.

In addition to the baryons the other particles which are involved in the nuclear interaction are called *mesons*. There are first the pions, which come in three varieties, positive, negative, and neutral; they form another multiplet. We have also found some new things called A'-mesons, and they occur as a doublet, K^+ and K^0 . Also, every particle has its antiparticle, unless a particle is its own antiparticle. For example, the π^- and the π^+ are antiparticles, but the π^0 is its own antiparticle. The K^- and K^+ are antiparticles, and the K^0 and \bar{K}^0 . In addition, in 1961 we also found some more mesons or *maybe* mesons which disintegrate almost immediately.

Table 2-2. Elementary Particles

MASS in Gev.	CHARGE -e	CHARGE 0	CHARGE +e	GROUP STRANG
1.4	Σ^-	Σ^0	Σ^+	$S = 2$
1.3	Ξ^-	Ξ^0		$S = 2$
1.2	Ξ^-	Ξ^0	Ξ^+	$S = 1$
1.1		Λ^0		$S = 1$
1.0				
0.9		n	p	$S = 0$
0.8		ω		$S = 0$
0.7	ρ^-	ρ^0	ρ^+	$S = 0$
0.6				
0.5	K^-	K^0	K^+	$S = -1$
0.4				
0.3				
0.2				
0.1	π^-	π^0	π^+	$S = 0$
0	e^-	ν		

A thing called w which goes into three pions has a mass 780 on this scale, and somewhat less certain is an object which disintegrates into two pions. These particles, called mesons and baryons, and the antiparticles of the mesons are on the same chart, but the antiparticles of the baryons must be put on another chart, "reflected" through the charge-zero column.

Just as Mendeleev's chart was very good, except for the fact that there were a number of rare earth elements which were hanging out loose from it, so we have a number of things hanging out loose from this chart—particles which do not interact strongly in nuclei, have nothing to do with a nuclear interaction, and do not have a strong interaction (I mean the powerful kind of interaction of nuclear energy). These are called leptons, and they are the following: there is the electron, which has a very small mass on this scale, only 0.510 Mev. Then there is that other, the μ -meson, the muon, which has a mass much higher, 206 times as heavy as an electron. So far as we can tell, by all experiments so far, the difference between the electron and the muon is nothing but the mass. Everything works exactly the same for the muon as for the electron, except that one is heavier than the other. Why is there another one heavier; what is the use for it? We do not know. In addition, there is a lepton which is neutral, called a neutrino, and this particle has zero mass. In fact, it is now known that there are *two* different kinds of neutrinos, one related to electrons and the other related to muons.

Finally, we have two other particles which do not interact strongly with the nuclear ones: one is a photon, and perhaps, if the field of gravity also has a quantum-mechanical analog (a quantum theory of gravitation has not yet been worked out), then there will be a particle, a graviton, which will have zero mass.

What is this "zero mass"? The masses given here are the masses of the particles *at rest*. The fact that a particle has zero mass means, in a way, that it cannot *be at rest*. A photon is never at rest, it is always moving at 186,000 miles a second. We will understand more what mass means when we understand the theory of relativity, which will come in due time.

Thus we are confronted with a large number of particles, which together seem to be the fundamental constituents of matter. Fortunately, these particles are not *all* different in their *interactions* with one another. In fact, there seem to be just *four kinds* of interaction between particles which, in the order of decreasing strength, are the nuclear force, electrical interactions, the beta-decay interaction, and gravity. The photon is coupled to all charged particles and the strength of the interaction is measured by some number, which is $1/137$. The detailed law of this coupling is known, that is quantum electrodynamics. • Gravity is coupled to all *energy*, but its coupling is extremely weak, much weaker than that of electricity. This law is also known. Then there are the so-called weak decays—beta decay, which causes the neutron to disintegrate into proton, electron, and neutrino, relatively slowly. This law is only partly known. The so-called strong interaction, the meson-baryon interaction, has a strength of 1 in this scale, and the law is completely unknown, although there are a number of known rules, such as that the number of baryons does not change in any reaction.

Table 2-3. Elementary Interactions

Coupling	Strength*	Law
Photon to charged particles	$\sim 10^{-2}$	Law known
Gravity to all energy	$\sim 10^{-40}$	Law known
Weak decays	$\sim 10^{-5}$	Law partly known
Mesons to baryons	~ 1	Law unknown (some rules known)

* The "strength" is a dimensionless measure of the coupling constant involved in each interaction (\sim means "approximately").

This then, is the horrible condition of our physics today. To summarize it, I would say this: outside the nucleus, we seem to know all; inside it, quantum mechanics is valid—the principles of quantum mechanics have not been found to fail. The stage on which we put all of our knowledge, we would say, is relativistic space-time; perhaps gravity is involved in space-time. We do not know how the universe got started, and we have never made experiments which check our ideas of space and time accurately, below some tiny distance, so we only *know* that our ideas work above that distance. We should also add that the rules of the game are the quantum-mechanical principles, and those principles apply, so far as we can tell, to the new particles as well as to the old. The origin of the forces in nuclei leads us to new particles, but unfortunately they appear in great profusion and we lack a complete understanding of their interrelationship, although we already know that there are some very surprising relationships among them. We seem gradually to be groping toward an understanding of the world of sub-atomic particles, but we really do not know how far we have yet to go in this task.

The Relation of Physics to Other Sciences

3-1 Introduction

Physics is the most fundamental and all-inclusive of the sciences, and has had a profound effect on all scientific development. In fact, physics is the present-day equivalent of what used to be called *natural philosophy*, from which most of our modern sciences arose. Students of many fields find themselves studying physics because of the basic role it plays in all phenomena. In this chapter we shall try to explain what the fundamental problems in the other sciences are, but of course it is impossible in so small a space really to deal with the complex, subtle, beautiful matters in these other fields. Lack of space also prevents our discussing the relation of physics to engineering, industry, society, and war, or even the most remarkable relationship between mathematics and physics. (Mathematics is not a science from our point of view, in the sense that it is not a *natural* science. The test of its validity is not experiment.) We must, incidentally, make it clear from the beginning that if a thing is not a science, it is not necessarily bad. For example, love is not a science. So, if something is said not to be a science, it does not mean that there is something wrong with it; it just means that it is not a science.

3-2 Chemistry

The science which is perhaps the most deeply affected by physics is chemistry. Historically, the early days of chemistry dealt almost entirely with what we now call inorganic chemistry, the chemistry of substances which are not associated with living things. Considerable analysis was required to discover the existence of the many elements and their relationships—how they make the various relatively simple compounds found in rocks, earth, etc. This early chemistry was very important for physics. The interaction between the two sciences was very great because the theory of atoms was substantiated to a large extent by experiments in chemistry. The theory of chemistry, i.e., of the reactions themselves, was summarized to a large extent in the periodic chart of Mendeleev, which brings out many strange relationships among the various elements, and it was the collection of rules as to which substance is combined with which, and how, that constituted inorganic chemistry. All these rules were ultimately explained in principle by quantum mechanics, so that theoretical chemistry is in fact physics. On the other hand, it must be emphasized that this explanation is *in principle*. We have already discussed the difference between knowing the rules of the game of chess, and being able to play. So it is that we may know the rules, but we cannot play very well. It turns out to be very difficult to predict precisely what will happen in a given chemical reaction; nevertheless, the deepest part of theoretical chemistry must end up in quantum mechanics.

There is also a branch of physics and chemistry which was developed by both sciences together, and which is extremely important. This is the method of statistics applied in a situation in which there are mechanical laws, which is aptly called *statistical mechanics*. In any chemical situation a large number of atoms are involved, and we have seen that the atoms are all jiggling around in a very random and complicated way. If we could analyze each collision, and be able to follow in detail the motion of each molecule, we might hope to figure out what would happen, but the many numbers needed to keep track of all these molecules exceeds so enormously the capacity of any computer, and certainly the capacity of

3-1 Introduction

3-2 Chemistry

3-3 Biology

3-4 Astronomy

3-5 Geology

3-6 Psychology

3-7 How did it get that way?

the mind, that it was important to develop a method for dealing with such complicated situations. Statistical mechanics, then, is the science of the phenomena of heat, or thermodynamics. Inorganic chemistry is, as a science, now reduced essentially to what are called physical chemistry and quantum chemistry; physical chemistry to study the rates at which reactions occur and what is happening in detail (How do the molecules hit? Which pieces fly off first?, etc.), and quantum chemistry to help us understand what happens in terms of the physical laws.

The other branch of chemistry is *organic chemistry*, the chemistry of the substances which are associated with living things. For a time it was believed that the substances which are associated with living things were so marvelous that they could not be made by hand, from inorganic materials. This is not at all true—they are just the same as the substances made in inorganic chemistry, but more complicated arrangements of atoms are involved. Organic chemistry obviously has a very close relationship to the biology which supplies its substances, and to industry, and furthermore, much physical chemistry and quantum mechanics can be applied to organic as well as to inorganic compounds. However, the main problems of organic chemistry are not in these aspects, but rather in the analysis and synthesis of the substances which are formed in biological systems, in living things. This leads imperceptibly, in steps, toward biochemistry, and then into biology itself, or molecular biology.

3-3 Biology

Thus we come to the science of *biology*, which is the study of living things. In the early days of biology, the biologists had to deal with the purely descriptive problem of finding out *what* living things there were, and so they just had to count such things as the hairs of the limbs of fleas. After these matters were worked out with a great deal of interest, the biologists went into the *machinery* inside the living bodies, first from a gross standpoint, naturally, because it takes some effort to get into the finer details.

There was an interesting early relationship between physics and biology in which biology helped physics in the discovery of the *conservation of energy*, which was first demonstrated by Mayer in connection with the amount of heat taken in and given out by a living creature.

If we look at the processes of biology of living animals more closely, we see *many* physical phenomena: the circulation of blood, pumps, pressure, etc. There are nerves: we know what is happening when we step on a sharp stone, and that somehow or other the information goes from the leg up. It is interesting how that happens. In their study of nerves, the biologists have come to the conclusion that nerves are very fine tubes with a complex wall which is very thin; through this wall the cell pumps ions, so that there are positive ions on the outside and negative ions on the inside, like a capacitor. Now this membrane has an interesting property; if it "discharges" in one place, i.e., if some of the ions were able to move through one place, so that the electric voltage is reduced there, that electrical influence makes itself felt on the ions in the neighborhood, and it affects the membrane in such a way that it lets the ions through at neighboring points also. This in turn affects it farther along, etc., and so there is a wave of "penetrability" of the membrane which runs down the fiber when it is "excited" at one end by stepping on the sharp stone. This wave is somewhat analogous to a long sequence of vertical dominoes; if the end one is pushed over, that one pushes the next, etc. Of course this will transmit only one message unless the dominoes are set up again; and similarly in the nerve cell, there are processes which pump the ions slowly out again, to get the nerve ready for the next impulse. So it is that we know what we are doing (or at least where we are). Of course the electrical effects associated with this nerve impulse can be picked up with electrical instruments, and because there *are* electrical effects, obviously the physics of electrical effects has had a great deal of influence on understanding the phenomenon.

The opposite effect is that, from somewhere in the brain, a message is sent out along a nerve. What happens at the end of the nerve? There the nerve branches

out into fine little things, connected to a structure near a muscle, called an end-plate. For reasons which are not exactly understood, when the impulse reaches the end of the nerve, little packets of a chemical called acetylcholine are shot off (five or ten molecules at a time) and they affect the muscle fiber and make it contract—how simple! What makes a muscle contract? A muscle is a very large number of fibers close together, containing two different substances, myosin and actomyosin, but the machinery by which the chemical reaction induced by acetylcholine can modify the dimensions of the molecule is not yet known. Thus the fundamental processes in the muscle that make mechanical motions are not known.

Biology is such an enormously wide field that there are hosts of other problems that we cannot mention at all—problems on how vision works (what the light does in the eye), how hearing works, etc. (The way in which *thinking* works we shall discuss later under psychology.) Now, these things concerning biology which we have just discussed are, from a biological standpoint, really not fundamental, at the bottom of life, in the sense that even if we understood them we still would not understand life itself. To illustrate: the men who study nerves feel their work is very important, because after all you cannot have animals without nerves. But you *can* have *life* without nerves. Plants have neither nerves nor muscles, but they are working, they are alive, just the same. So for the fundamental problems of biology we must look deeper; when we do, we discover that all living things have a great many characteristics in common. The most common feature is that they are made of *cells*, within each of which is complex machinery for doing things chemically. In plant cells, for example, there is machinery for picking up light and generating sucrose, which is consumed in the dark to keep the plant alive. When the plant is eaten the sucrose itself generates in the animal a series of chemical reactions very closely related to photosynthesis (and its opposite effect in the dark) in plants.

In the cells of living systems there are many elaborate chemical reactions, in which one compound is changed into another and another. To give some impression of the enormous efforts that have gone into the study of biochemistry, the chart in Fig. 3-1 summarizes our knowledge to date on just one small part of the many series of reactions which occur in cells, perhaps a percent or so of it.

Here we see a whole series of molecules which change from one to another in a sequence or cycle of rather small steps. It is called the Krebs cycle, the respiratory cycle. Each of the chemicals and each of the steps is fairly simple, in terms of what change is made in the molecule, but—and this is a centrally important discovery in biochemistry—these changes are *relatively difficult to accomplish in a laboratory*. If we have one substance and another very similar substance, the one does not just turn into the other, because the two forms are usually separated by

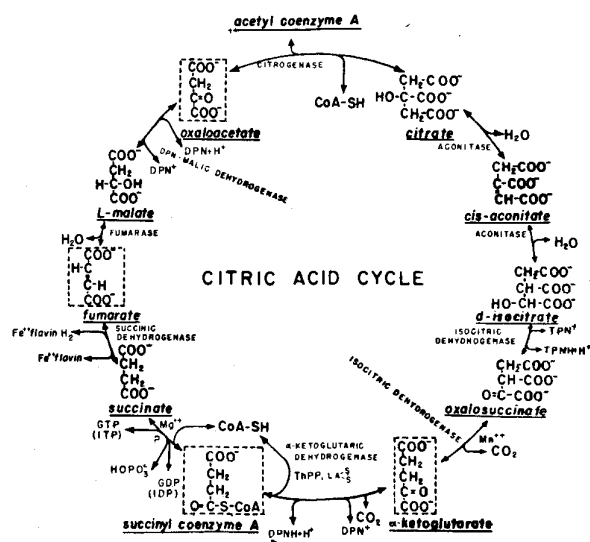


Fig. 3-1. The Krebs cycle.

an energy barrier or "hill." Consider this analogy: If we wanted to take an object from one place to another, at the same level but on the other side of a hill, we could push it over the top, but to do so requires the addition of some energy. Thus most chemical reactions do not occur, because there is what is called an *activation energy* in the way. In order to add an extra atom to our chemical requires that we get it *close* enough that some rearrangement can occur; then it will stick. But if we cannot give it enough energy to get it close enough, it will not go to completion, it will just go part way up the "hill" and back down again. However, if we could literally take the molecules in our hands and push and pull the atoms around in such a way as to open a hole to let the new atom in, and then let it snap back, we would have found another way, *around* the hill, which would not require extra energy, and the reaction would go easily. Now there actually *are*, in the cells, *very* large molecules, much larger than the ones whose changes we have been describing, which in some complicated way hold the smaller molecules just right, so that the reaction can occur easily. These very large and complicated things are called *enzymes*. (They were first called ferments, because they were originally discovered in the fermentation of sugar. In fact, some of the first reactions in the cycle were discovered there.) In the presence of an enzyme the reaction will go.

An enzyme is made of another substance called *protein*. Enzymes are very big and complicated, and each one is different, each being built to control a certain special reaction. The names of the enzymes are written in Fig. 3-1 at each reaction. (Sometimes the same enzyme may control two reactions.) We emphasize that the enzymes themselves are not involved in the reaction directly. They do not change; they merely let an atom go from one place to another. Having done so, the enzyme is ready to do it to the next molecule, like a machine in a factory. Of course, there must be a supply of certain atoms and a way of disposing of other atoms. Take hydrogen, for example: there are enzymes which have special units on them which carry the hydrogen for all chemical reactions. For example, there are three or four hydrogen-reducing enzymes which are used all over our cycle in different places. It is interesting that the machinery which liberates some hydrogen at one place will take that hydrogen and use it somewhere else.

The most important feature of the cycle of Fig. 3-1 is the transformation from GDP to GTP (guanadine-di-phosphate to guanadine-tri-phosphate) because the one substance has much more energy in it than the other. Just as there is a "box" in certain enzymes for carrying hydrogen atoms around, there are special *energy-carrying* "boxes" which involve the triphosphate group. So, GTP has more energy than GDP and if the cycle is going one way, we are producing molecules which have extra energy and which can go drive some other cycle which *requires* energy, for example the contraction of muscle. The muscle will not contract unless there is GTP. We can take muscle fiber, put it in water, and add GTP, and the fibers contract, changing GTP to GDP if the right enzymes are present. So the real system is in the GDP-GTP transformation; in the dark the GTP which has been stored up during the day is used to run the whole cycle around the other way. An enzyme you see, does not care in which direction the reaction goes, for if it did it would violate one of the laws of physics.

Physics is of great importance in biology and other sciences for still another reason, that has to do with *experimental techniques*. In fact, if it were not for the great development of experimental physics, these biochemistry charts would not be known today. The reason is that the most useful tool of all for analyzing this fantastically complex system is to *label* the atoms which are used in the reactions. Thus, if we could introduce into the cycle some carbon dioxide which has a "green mark" on it, and then measure after three seconds where the green mark is, and again measure after ten seconds, etc., we could trace out the course of the reactions. What are the "green marks"? They are different *isotopes*. We recall that the chemical properties of atoms are determined by the number of *electrons*, not by the mass of the nucleus. But there can be, for example in carbon, six neutrons or seven neutrons, together with the six protons which all carbon nuclei have. Chemically, the two atoms C_{12} and C_{13} are the same, but they differ in weight and they have different nuclear properties, and so they are distinguishable.

By using these isotopes of different weights, or even radioactive isotopes like C^{14} , which provide a more sensitive means for tracing very small quantities, it is possible to trace the reactions.

Now, we return to the description of enzymes and proteins. All proteins are not enzymes, but all enzymes are proteins. There are many proteins, such as the proteins in muscle, the structural proteins which are, for example, in cartilage and hair, skin, etc., that are not themselves enzymes. However, proteins are a very characteristic substance of life: first of all they make up all the enzymes, and second, they make up much of the rest of living material. Proteins have a very interesting and simple structure. They are a series, or chain, of different *amino acids*. There are twenty different amino acids, and they all can combine with each other to form chains in which the backbone is $CO-NH$, etc. Proteins are nothing but chains of various ones of these twenty amino acids. Each of the amino acids probably serves some special purpose. Some, for example, have a sulphur atom at a certain place; when two sulphur atoms are in the same protein, they form a bond, that is, they tie the chain together at two points and form a loop. Another has extra oxygen atoms which make it an acidic substance, another has a basic characteristic. Some of them have big groups hanging out to one side, so that they take up a lot of space. One of the amino acids, called proline, is not really an amino acid, but imino acid. There is a slight difference, with the result that when proline is in the chain, there is a kink in the chain. If we wished to manufacture a particular protein, we would give these instructions: put one of those sulphur hooks here; next, add something to take up space; then attach something to put a kink in the chain. In this way, we will get a complicated-looking chain, hooked together and having some complex structure; this is presumably just the manner in which all the various enzymes are made. One of the great triumphs in recent times (since 1960), was at last to discover the exact spatial atomic arrangement of certain proteins, which involve some fifty-six or sixty amino acids in a row. Over a thousand atoms (more nearly two thousand, if we count the hydrogen atoms) have been located in a complex pattern in two proteins. The first was hemoglobin. One of the sad aspects of this discovery is that we cannot see anything from the pattern; we do not understand why it works the way it does. Of course, that is the next problem to be attacked.

Another problem is how do the enzymes know what to be? A red-eyed fly makes a red-eyed fly baby, and so the information for the whole pattern of enzymes to make red pigment must be passed from one fly to the next. This is done by a substance in the nucleus of the cell, not a protein, called DNA (short for deoxyribose nucleic acid). This is the key substance which is passed from one cell to another (for instance, sperm cells consist mostly of DNA) and carries the information as to how to make the enzymes. DNA is the "blueprint." What does the blueprint look like and how does it work? First, the blueprint must be able to reproduce itself. Secondly, it must be able to instruct the protein. Concerning the reproduction, we might think that this proceeds like cell reproduction. Cells simply grow bigger and then divide in half. Must it be thus with DNA molecules, then, that they too grow bigger and divide in half? Every *atom* certainly does not grow bigger and divide in half! No, it is impossible to reproduce a molecule except by some more clever way.

The structure of the substance DNA was studied for a long time, first chemically to find the composition, and then with x-rays to find the pattern in space. The result was the following remarkable discovery: The DNA molecule is a pair of chains, twisted upon each other. The backbone of each of these chains, which are analogous to the chains of proteins but chemically quite different, is a series of sugar and phosphate groups, as shown in Fig. 3-2. Now we see how the chain can contain instructions, for if we could split this chain down the middle, we would have a series *BAADC* . . . and every living thing could have a different series. Thus perhaps, in some way, the specific *instructions* for the manufacture of proteins are contained in the specific *series* of the DNA.

Attached to each sugar along the line, and linking the two chains together, are certain pairs of cross-links. However, they are not all of the same kind; there are

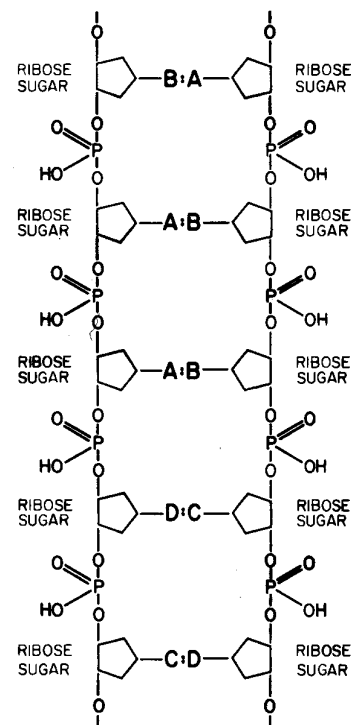


Fig. 3-2. Schematic diagram of DNA.

four kinds, called adenine, thymine, cytosine, and guanine, but let us call them *A*, *B*, *C*, and *D*. The interesting thing is that only certain pairs can sit opposite each other, for example *A* with *B* and *C* with *D*. These pairs are put on the two chains in such a way that they "fit together," and have a strong energy of interaction. However, *C* will not fit with *A*, and *B* will not fit with *C*; they will only fit in pairs, *A* against *B* and *C* against *D*. Therefore if one is *C*, the other must be *D*, etc. Whatever the letters may be in one chain, each one must have its specific complementary letter on the other chain.

What then about reproduction? Suppose we split this chain in two. How can we make another one just like it? If, in the substances of the cells, there is a manufacturing department which brings up phosphate, sugar, and *A*, *B*, *C*, *D* units not connected in a chain, the only ones which will attach to our split chain will be the correct ones, the complements of *BAADC* . . . , namely, *ABBCD* . . . Thus what happens is that the chain splits down the middle during cell division, one half ultimately to go with one cell, the other half to end up in the other cell; when separated, a new complementary chain is made by each half-chain.

Next comes the question, precisely how does the order of the *A*, *B*, *C*, *D* units determine the arrangement of the amino acids in the protein? This is the central unsolved problem in biology today. The first clues, or pieces of information, however, are these: There are in the cell tiny particles called microsomes, and it is now known that that is the place where proteins are made. But the microsomes are not in the nucleus, where the DNA and its instructions are. Something seems to be the matter. However, it is also known that little molecule pieces come off the DNA—not as long as the big DNA molecule that carries all the information itself, but like a small section of it. This is called RNA, but that is not essential. It is a kind of copy of the DNA, a short copy. The RNA, which somehow carries a message as to what kind of protein to make goes over to the microsome; that is known. When it gets there, protein is synthesized at the microsome. That is also known. However, the details of how the amino acids come in and are arranged in accordance with a code that is on the RNA are, as yet, still unknown. We do not know how to read it. If we knew, for example, the "lineup" *A*, *B*, *C*, *C*, *A*, we could not tell you what protein is to be made.

Certainly no subject or field is making more progress on so many fronts at the present moment, than biology, and if we were to name the most powerful assumption of all, which leads one on and on in an attempt to understand life, it is that *all things are made of atoms*, and that everything that living things do can be understood in terms of the jiggings and wiggings of atoms.

3-4 Astronomy

In this rapid-fire explanation of the whole world, we must now turn to astronomy. Astronomy is older than physics. In fact, it got physics started by showing the beautiful simplicity of the motion of the stars and planets, the understanding of which was the *beginning* of physics. But the most remarkable discovery in all of astronomy is that *the stars are made of atoms of the same kind as those on the earth** How was this done? Atoms liberate light which has definite fre-

* How I'm rushing through this! How much each sentence in this brief story contains. "The stars are made of the same atoms as the earth." I usually pick one small topic like this to give a lecture on. Poets say science takes away from the beauty of the stars—mere globs of gas atoms. Nothing is "mere." I too can see the stars on a desert night, and feel them. But do I see less or more? The vastness of the heavens stretches my imagination—stuck on this carousel my little eye can catch one-million-year-old light. A vast pattern—of which I am a part—perhaps my stuff was belched from some forgotten star, as one is belching there. Or see them with the greater eye of Palomar, rushing all apart from some common starting point when they were perhaps all together. What is the pattern, or the meaning, or the *why*? It does not do harm to the mystery to know a little about it. For far more marvelous is the truth than any artists of the past imagined! Why do the poets of the present not speak of it? What men are poets who can speak of Jupiter if he were like a man, but if he is an immense spinning sphere of methane and ammonia must be silent?

quencies, something like the timbre of a musical instrument, which has definite pitches or frequencies of sound. When we are listening to several different tones we can tell them apart, but when we look with our eyes at a mixture of colors we cannot tell the parts from which it was made, because the eye is nowhere near as discerning as the ear in this connection. However, with a spectroscope we *can* analyze the frequencies of the light waves and in this way we can see the very tunes of the atoms that are in the different stars. As a matter of fact, two of the chemical elements were discovered on a star before they were discovered on the earth. Helium was discovered on the sun, whence its name, and technetium was discovered in certain cool stars. This, of course, permits us to make headway in understanding the stars, because they are made of the same kinds of atoms which are on the earth. Now we know a great deal about the atoms, especially concerning their behavior under conditions of high temperature but not very great density, so that we can analyze by statistical mechanics the behavior of the stellar substance. Even though we cannot reproduce the conditions on the earth, using the basic physical laws we often can tell precisely, or very closely, what will happen. So it is that physics aids astronomy. Strange as it may seem, we understand the distribution of matter in the interior of the sun far better than we understand the interior of the earth. What goes on *inside* a star is better understood than one might guess from the difficulty of having to look at a little dot of light through a telescope, because we can *calculate* what the atoms in the stars should do in most circumstances.

One of the most impressive discoveries was the origin of the energy of the stars, that makes them continue to burn. One of the men who discovered this was out with his girl friend the night after he realized that *nuclear reactions* must be going on in the stars in order to make them shine. She said "Look at how pretty the stars shine!" He said "Yes, and right now I am the only man in the world who knows *why* they shine." She merely laughed at him. She was not impressed with being out with the only man who, at that moment, knew why stars shine. Well, it is sad to be alone, but that is the way it is in this world.

It is the nuclear "burning" of hydrogen which supplies the energy of the sun; the hydrogen is converted into helium. Furthermore, ultimately, the manufacture of various chemical elements proceeds in the centers of the stars, from hydrogen. The stuff of which *we* are made, was "cooked" once, in a star, and spit out. How do we know? Because there is a clue. The proportion of the different isotopes—how much C^{12} , how much C^{13} , etc., is something which is never changed by *chemical* reactions, because the chemical reactions are so much the same for the two. The proportions are purely the result of *nuclear* reactions. By looking at the proportions of the isotopes in the cold, dead ember which we are, we can discover what the *furnace* was like in which the stuff of which we are made was formed. That furnace was like the stars, and so it is very likely that our elements were "made" in the stars and spit out in the explosions which we call novae and supernovae. Astronomy is so close to physics that we shall study many astronomical things as we go along.

3-5 Geology

We turn now to what are called *earth sciences*, or *geology*. First, meteorology and the weather. Of course the *instruments* of meteorology are physical instruments, and the development of experimental physics made these instruments possible, as was explained before. However, the theory of meteorology has never been satisfactorily worked out by the physicist. "Well," you say, "there is nothing but air, and we know the equations of the motions of air." Yes we do. "So if we know the condition of air today, why can't we figure out the condition of the air tomorrow?" First, we do not *really* know what the condition is today, because the air is swirling and twisting everywhere. It turns out to be very sensitive, and even unstable. If you have ever seen water run smoothly over a dam, and then turn into a large number of blobs and drops as it falls, you will understand what I mean by unstable. You know the condition of the water before it goes over the

spillway; it is perfectly smooth; but the moment it begins to fall, where do the drops begin? What determines how big the lumps are going to be and where they will be? That is not known, because the water is unstable. Even a smooth moving mass of air, in going over a mountain turns into complex whirlpools and eddies. In many fields we find this situation of *turbulent flow* that we cannot analyze today. Quickly we leave the subject of weather, and discuss geology!

The question basic to geology is, what makes the earth the way it is? The most obvious processes are in front of your very eyes, the erosion processes of the rivers, the winds, etc. It is easy enough to understand these, but for every bit of erosion there is an equal amount of something else going on. Mountains are no lower today, on the average, than they were in the past. There must be *mountain-forming* processes. You will find, if you study geology, that there *are* mountain-forming processes and vulcanism, which nobody understands but which is half of geology. The phenomenon of volcanoes is really not understood. What makes an earthquake is, ultimately, not understood. It is understood that if something is pushing something else, it snaps and will slide—that is all right. But what pushes, and why? The theory is that there are currents inside the earth—circulating currents, due to the difference in temperature inside and outside—which, in their motion, push the surface slightly. Thus if there are two opposite circulations next to each other, the matter will collect in the region where they meet and make belts of mountains which are in unhappy stressed conditions, and so produce volcanoes and earthquakes.

What about the inside of the earth? A great deal is known about the speed of earthquake waves through the earth and the density of distribution of the earth. However, physicists have been unable to get a good theory as to how dense a substance should be at the pressures that would be expected at the center of the earth. In other words, we cannot figure out the properties of matter very well in these circumstances. We do much less well with the earth than we do with the conditions of matter in the stars. The mathematics involved seems a little too difficult, so far, but perhaps it will not be too long before someone realizes that it is an important problem, and really work it out. The other aspect, of course, is that even if we did know the density, we cannot figure out the circulating currents. Nor can we really work out the properties of rocks at high pressure. We cannot tell how fast the rocks should "give"; that must all be worked out by experiment.

3-6 Psychology

Next, we consider the science of *psychology*. Incidentally, psychoanalysis is not a science: it is at best a medical process, and perhaps even more like witch-doctoring. It has a theory as to what causes disease—lots of different "spirits," etc. The witch doctor has a theory that a disease like malaria is caused by a spirit which comes into the air; it is not cured by shaking a snake over it, but quinine does help malaria. So, if you are sick, I would advise that you go to the witch doctor because he is the man in the tribe who knows the most about the disease; on the other hand, his knowledge is not science. Psychoanalysis has not been checked carefully by experiment, and there is no way to find a list of the number of cases in which it works, the number of cases in which it does not work, etc.

The other branches of psychology, which involve things like the physiology of sensation—what happens in the eye, and what happens in the brain—are, if you wish, less interesting. But some small but real progress has been made in studying them. One of the most interesting technical problems may or may not be called psychology. The central problem of the mind, if you will, or the nervous system, is this: when an animal learns something, it can do something different than it could before, and its brain cell must have changed too, if it is made out of atoms. *In what way is it different?* We do not know where to look, or what to look for, when something is memorized. We do not know what it means, or what change there is in the nervous system, when a fact is learned. This is a very important problem which has not been solved at all. Assuming, however, that there is some kind of memory thing, the brain is such an enormous mass of interconnect-

ing wires and nerves that it probably cannot be analyzed in a straightforward manner. There is an analog of this to computing machines and computing elements, in that they also have a lot of lines, and they have some kind of element, analogous, perhaps, to the synapse, or connection of one nerve to another. This is a very interesting subject which we have not the time to discuss further—the relationship between thinking and computing machines. It must be appreciated, of course, that this subject will tell us very little about the real complexities of ordinary human behavior. All human beings are so different. It will be a long time before we get there. We must start much further back. If we could even figure out how a *dog* works, we would have gone pretty far. Dogs are easier to understand, but nobody yet knows how dogs work.

3-7 How did it get that way?

In order for physics to be useful to other sciences in a *theoretical* way, other than in the invention of instruments, the science in question must supply to the physicist a description of the object in a physicist's language. They can say "why does a frog jump?," and the physicist cannot answer. If they tell him what a frog is, that there are so many molecules, there is a nerve here, etc., that is different. If they will tell us, more or less, what the earth or the stars are like, then we can figure it out. In order for physical theory to be of any use, we must know where the atoms are located. In order to understand the chemistry, we must know exactly what atoms are present, for otherwise we cannot analyze it. That is but one limitation, of course.

There is another *kind* of problem in the sister sciences which does not exist in physics; we might call it, for lack of a better term, the historical question. How did it get that way? If we understand all about biology, we will want to know how all the things which are on the earth got there. There is the theory of evolution, an important part of biology. In geology, we not only want to know how the mountains are forming, but how the entire earth was formed in the beginning, the origin of the solar system, etc. That, of course, leads us to want to know what kind of matter there was in the world. How did the stars evolve? What were the initial conditions? That is the problem of astronomical history. A great deal has been found out about the formation of stars, the formation of elements from which we were made, and even a little about the origin of the universe.

There is no historical question being studied in physics at the present time. We do not have a question, "Here are the laws of physics, how did they get that way?" We do not imagine, at the moment, that the laws of physics are somehow changing with time, that they were different in the past than they are at present. Of course they *may* be, and the moment we find they *are*, the historical question of physics will be wrapped up with the rest of the history of the universe, and then the physicist will be talking about the same problems as astronomers, geologists, and biologists.

Finally, there is a physical problem that is common to many fields, that is very old, and that has not been solved. It is not the problem of finding new fundamental particles, but something left over from a long time ago—over a hundred years. Nobody in physics has really been able to analyze it mathematically satisfactorily in spite of its importance to the sister sciences. It is the analysis of *circulating or turbulent fluids*. If we watch the evolution of a star, there comes a point where we can deduce that it is going to start convection, and thereafter we can no longer deduce what should happen. A few million years later the star explodes, but we cannot figure out the reason. We cannot analyze the weather. We do not know the patterns of motions that there should be inside the earth. The simplest form of the problem is to take a pipe that is very long and push water through it at high speed. We ask: to push a given amount of water through that pipe, how much pressure is needed? No one can analyze it from first principles and the properties of water. If the water flows very slowly, or if we use a thick goo like honey, then we can do it nicely. You will find that in your textbook.

What we really cannot do is deal with actual, wet water running through a pipe. That is the central problem which we ought to solve some day, and we have not.

A poet once said, "The whole universe is in a glass of wine." We will probably never know in what sense he meant that, for poets do not write to be understood. But it is true that if we look at a glass of wine closely enough we see the entire universe. There are the things of physics: the twisting liquid which evaporates depending on the wind and weather, the reflections in the glass, and our imagination adds the atoms. The glass is a distillation of the earth's rocks, and in its composition we see the secrets of the universe's age, and the evolution of stars. What strange array of chemicals are in the wine? How did they come to be? There are the ferments, the enzymes, the substrates, and the products. There in wine is found the great generalization: all life is fermentation. Nobody can discover the chemistry of wine without discovering, as did Louis Pasteur, the cause of much disease. How vivid is the claret, pressing its existence into the consciousness that watches it! If our small minds, for some convenience, divide this glass of wine, this universe, into parts—physics, biology, geology, astronomy, psychology, and so on—remember that nature does not know it! So let us put it all back together, not forgetting ultimately what it is for. Let it give us one more final pleasure: drink it and forget it all!

Conservation of Energy

4-1 What is energy?

In this chapter, we begin our more detailed study of the different aspects of physics, having finished our description of things in general. To illustrate the ideas and the kind of reasoning that might be used in theoretical physics, we shall now examine one of the most basic laws of physics, the conservation of energy.

There is a fact, or if you wish, a *law*, governing all natural phenomena that are known to date. There is no known exception to this law—it is exact so far as we know. The law is called the *conservation of energy*. It states that there is a certain quantity, which we call energy, that does not change in the manifold changes which nature undergoes. That is a most abstract idea, because it is a mathematical principle; it says that there is a numerical quantity which does not change when something happens. It is not a description of a mechanism, or anything concrete; it is just a strange fact that we can calculate some number and when we finish watching nature go through her tricks and calculate the number again, it is the same. (Something like the bishop on a red square, and after a number of moves—details unknown—it is still on some red square. It is a law of this nature.) Since it is an abstract idea, we shall illustrate the meaning of it by an analogy.

Imagine a child, perhaps "Dennis the Menace," who has blocks which are absolutely indestructible, and cannot be divided into pieces. Each is the same as the other. Let us suppose that he has 28 blocks. His mother puts him with his 28 blocks into a room at the beginning of the day. At the end of the day, being curious, she counts the blocks very carefully, and discovers a phenomenal law—no matter what he does with the blocks, there are always 28 remaining! This continues for a number of days, until one day there are only 27 blocks, but a little investigating shows that there is one under the rug—she must look everywhere to be sure that the number of blocks has not changed. One day, however, the number appears to change—there are only 26 blocks. Careful investigation indicates that the window was open, and upon looking outside, the other two blocks are found. Another day, careful count indicates that there are 30 blocks! This causes considerable consternation, until it is realized that Bruce came to visit, bringing his blocks with him, and he left a few at Dennis' house. After she has disposed of the extra blocks, she closes the window, does not let Bruce in, and then everything is going along all right, until one time she counts and finds only 25 blocks. However, there is a box in the room, a toy box, and the mother goes to open the toy box, but the boy says "No, do not open my toy box," and screams. Mother is not allowed to open the toy box. Being extremely curious, and somewhat ingenious, she invents a scheme! She knows that a block weighs three ounces, so she weighs the box at a time when she sees 28 blocks, and it weighs 16 ounces. The next time she wishes to check, she weighs the box again, subtracts sixteen ounces and divides by three. She discovers the following:

$$\left(\begin{array}{c} \text{number of} \\ \text{blocks seen} \end{array} \right) + \frac{(\text{weight of box}) - 16 \text{ ounces}}{3 \text{ ounces}} = \text{constant.} \quad (4.1)$$

There then appear to be some new deviations, but careful study indicates that the dirty water in the bathtub is changing its level. The child is throwing blocks into the water, and she cannot see them because it is so dirty, but she can find out how many blocks are in the water by adding another term to her formula. Since the original height of the water was 6 inches and each block raises the water a quarter

4-1 What is energy?

4-2 Gravitational potential energy

4-3 Kinetic energy

4-4 Other forms of energy

of an inch, this new formula would be:

$$\left(\frac{\text{number of}}{\text{blocks seen}} \right) + \frac{(\text{weight of box}) - 16 \text{ ounces}}{3 \text{ ounces}} + \frac{(\text{height of water}) - 6 \text{ inches}}{1/4 \text{ inch}} = \text{constant.} \quad (4.2)$$

In the gradual increase in the complexity of her world, she finds a whole series of terms representing ways of calculating how many blocks are in places where she is not allowed to look. As a result, she finds a complex formula, a quantity which *has to be computed*, which always stays the same in her situation.

What is the analogy of this to the conservation of energy? The most remarkable aspect that must be abstracted from this picture is that *there are no blocks*. Take away the first terms in (4.1) and (4.2) and we find ourselves calculating more or less abstract things. The analogy has the following points. First, when we are calculating the energy, sometimes some of it leaves the system and goes away, or sometimes some comes in. In order to verify the conservation of energy, we must be careful that we have not put any in or taken any out. Second, the energy has a large number of *different forms*, and there is a formula for each one. These are: gravitational energy, kinetic energy, heat energy, elastic energy, electrical energy, chemical energy, radiant energy, nuclear energy, mass energy. If we total up the formulas for each of these contributions, it will not change except for energy going in and out.

It is important to realize that in physics today, we have no knowledge of what energy *is*. We do not have a picture that energy comes in little blobs of a definite amount. It is not that way. However, there are formulas for calculating some numerical quantity, and when we add it all together it gives "28"—always the same number. It is an abstract thing in that it does not tell us the mechanism or the *reasons* for the various formulas.

4-2 Gravitational potential energy

Conservation of energy can be understood only if we have the formula for all of its forms. I wish to discuss the formula for gravitational energy near the surface of the Earth, and I wish to derive this formula in a way which has nothing to do with history but is simply a line of reasoning invented for this particular lecture to give you an illustration of the remarkable fact that a great deal about nature can be extracted from a few facts and close reasoning. It is an illustration of the kind of work theoretical physicists become involved in. It is patterned after a most excellent argument by Mr. Carnot on the efficiency of steam engines.*

Consider weight-lifting machines—machines which have the property that they lift one weight by lowering another. Let us also make a hypothesis: that *there is no such thing as perpetual motion* with these weight-lifting machines. (In fact, that there is no perpetual motion at all is a general statement of the law of conservation of energy.) We must be careful to define perpetual motion. First, let us do it for weight-lifting machines. If, when we have lifted and lowered a lot of weights and restored the machine to the original condition, we find that the net result is to have *lifted a weight*, then we have a perpetual motion machine because we can use that lifted weight to run something else. That is, *provided* the machine which lifted the weight is brought back to its exact *original condition*, and furthermore that it is completely *self-contained*—that it has not received the energy to lift that weight from some external source—like Bruce's blocks.

A very simple weight-lifting machine is shown in Fig. 4-1. This machine lifts weights three units "strong." We place three units on one balance pan, and one unit on the other. However, in order to get it actually to work, we must lift a little weight off the left pan. On the other hand, we could lift a one-unit weight



Fig. 4-1. Simple weight-lifting machine.

* Our point here is not so much the result, (4.3), which in fact you may already know, as the possibility of arriving at it by theoretical reasoning.

by lowering the three-unit weight, if we cheat a little by lifting a little weight off the other pan. Of course, we realize that with any *actual* lifting machine, we must add a little extra to get it to run. This we disregard, *temporarily*. Ideal machines, although they do not exist, do not require anything extra. A machine that we actually use can be, in a sense, *almost* reversible: that is, if it will lift the weight of three by lowering a weight of one, then it will also lift nearly the weight of one the same amount by lowering the weight of three.

We imagine that there are two classes of machines, those that are *not* reversible, which includes all real machines, and those that *are* reversible, which of course are actually not attainable no matter how careful we may be in our design of bearings, levers, etc. We suppose, however, that there is such a thing—a reversible machine—which lowers one unit of weight (a pound or any other unit) by one unit of distance, and at the same time lifts a three-unit weight. Call this reversible machine, Machine A. Suppose this particular reversible machine lifts the three-unit weight a distance X . Then suppose we have another machine, Machine B, which is not necessarily reversible, which also lowers a unit weight a unit distance, but which lifts three units a distance Y . We can now prove that Y is not higher than X ; that is, it is impossible to build a machine that will lift a weight *any higher* than it will be lifted by a reversible machine. Let us see why. Let us suppose that Y were higher than X . We take a one-unit weight and lower it one unit height with Machine B, and that lifts the three-unit weight up a distance Y . Then we could lower the weight from Y to X , obtaining *free power*, and use the reversible Machine A, running backwards, to lower the three-unit weight a distance X and lift the one-unit weight by one unit height. This will put the one-unit weight back where it was before, and leave both machines ready to be used again! We would therefore have perpetual motion if Y were higher than X , which we assumed was impossible. With those assumptions, we thus deduce that Y is not higher than X , so that of all machines that can be designed, the reversible machine is the best.

We can also see that all reversible machines must lift to *exactly the same height*. Suppose that B were really reversible also. The argument that Y is not higher than X is, of course, just as good as it was before, but we can also make our argument the other way around, using the machines in the opposite order, and prove that X is not higher than Y . This, then, is a very remarkable observation because it permits us to analyze the height to which different machines are going to lift something *without looking at the interior mechanism*. We know at once that if somebody makes an enormously elaborate series of levers that lift three units a certain distance by lowering one unit by one unit distance, and we compare it with a simple lever which does the same thing and is fundamentally reversible, his machine will lift it no higher, but perhaps less high. If his machine is reversible, we also know exactly *how high* it will lift. To summarize: every reversible machine, no matter how it operates, which drops one pound one foot and lifts a three-pound weight always lifts it the same distance, X . This is clearly a universal law of great utility. The next question is, of course, what is X ?

Suppose we have a reversible machine which is going to lift this distance X , three for one. We set up three balls in a rack which does not move, as shown in Fig. 4-2. One ball is held on a stage at a distance one foot above the ground. The machine can lift three balls, lowering one by a distance 1. Now, we have arranged that the platform which holds three balls has a floor and two shelves, exactly spaced at distance X , and further, that the rack which holds the balls is spaced at distance X , (a). First we roll the balls horizontally from the rack to the shelves, (b), and we suppose that this takes no energy because we do not change the height. The reversible machine then operates: it lowers the single ball to the floor, and it lifts the rack a distance X , (c). Now we have ingeniously arranged the rack so that these balls are again even with the platforms. Thus we unload the balls onto the rack, (d); having unloaded the balls, we can restore the machine to its original condition. Now we have three balls on the upper three shelves and one at the bottom. But the strange thing is that, in a certain way of speaking, we have not lifted *two* of them at all because, after all, there were balls on shelves 2 and 3

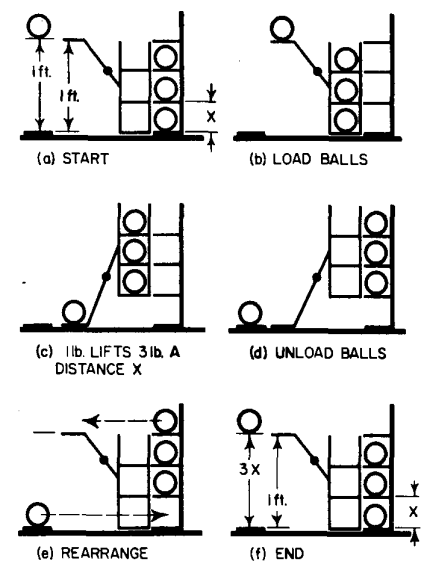


Fig. 4-2. A reversible machine.

before. The resulting effect has been to lift *one ball* a distance $3X$. Now, if $3X$ exceeds one foot, then we can *lower* the ball to return the machine to the initial condition, (f), and we can run the apparatus again. Therefore $3X$ *cannot* exceed one foot, for if $3X$ exceeds one foot we can make perpetual motion. Likewise, we can prove that *one foot cannot exceed $3X$* , by making the whole machine run the opposite way, since it is a reversible machine. Therefore $3X$ is neither *greater nor less than a foot*, and we discover then, by argument alone, the law that $X = \frac{1}{3}$ foot. The generalization is clear: one pound falls a certain distance in operating a reversible machine; then the machine can lift p pounds this distance divided by p . Another way of putting the result is that three pounds times the height lifted, which in our problem was X , is equal to one pound times the distance lowered, which is one foot in this case. If we take all the weights and multiply them by the heights at which they are now, above the floor, let the machine operate, and then multiply all the weights by all the heights again, *there will be no change*. (We have to generalize the example where we moved only one weight to the case where when we lower one we lift several different ones—but that is easy.)

We call the sum of the weights times the heights *gravitational potential energy*—the energy which an object has because of its relationship in space, relative to the earth. The formula for gravitational energy, then, so long as we are not too far from the earth (the force weakens as we go higher) is

$$\left(\begin{array}{c} \text{gravitational} \\ \text{potential energy} \\ \text{for one object} \end{array} \right) = (\text{weight}) \times (\text{height}). \quad (4.3)$$

It is a very beautiful line of reasoning. The only problem is that perhaps it is not true. (After all, nature does not *have* to go along with our reasoning.) For example, perhaps perpetual motion is, in fact, possible. Some of the assumptions may be wrong, or we may have made a mistake in reasoning, so it is always necessary to check. *It turns out experimentally*, in fact, to be true.

The general name of energy which has to do with location relative to something else is called *potential energy*. In this particular case, of course, we call it *gravitational potential energy*. If it is a question of electrical forces against which we are working, instead of gravitational forces, if we are "lifting" charges away from other charges with a lot of levers, then the energy content is called *electrical potential energy*. The general principle is that the change in the energy is the force times the distance that the force is pushed, and that this is a change in energy in general:

$$\left(\begin{array}{c} \text{change in} \\ \text{energy} \end{array} \right) = (\text{force}) \times \left(\begin{array}{c} \text{distance force} \\ \text{acts through} \end{array} \right). \quad (4.4)$$

We will return to many of these other kinds of energy as we continue the course.

The principle of the conservation of energy is very useful for deducing what will happen in a number of circumstances. In high school we learned a lot of laws about pulleys and levers used in different ways. We can now see that these "laws" are *all the same thing*, and that we did not have to memorize 75 rules to figure it out. A simple example is a smooth inclined plane which is, happily, a three-four-five triangle (Fig. 4-3). We hang a one-pound weight on the inclined plane with a pulley, and on the other side of the pulley, a weight W . We want to know how heavy W must be to balance the one pound on the plane. How can we figure that out? If we say it is just balanced, it is reversible and so can move up and down, and we can consider the following situation. In the initial circumstance, (a), the one pound weight is at the bottom and weight W is at the top. When W has slipped down in a reversible way, we have a one-pound weight at the top and the weight W the slant distance, (b), or five feet, from the plane in which it was before. We *lifted* the one-pound weight only *three* feet and we lowered W pounds by *five* feet. Therefore $W = \frac{3}{5}$ of a pound. Note that we deduced this from the *conservation of energy*, and not from force components. Cleverness, however, is relative. It can be deduced in a way which is even more brilliant, discovered by

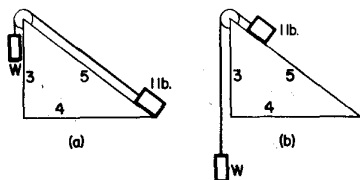


Fig. 4-3. Inclined plane.

Stevinus and inscribed on his tombstone. Figure 4-4 explains that it has to be $\frac{2}{3}$ of a pound, because the chain does not go around. It is evident that the lower part of the chain is balanced by itself, so that the pull of the five weights on one side must balance the pull of three weights on the other, or whatever the ratio of the legs. You see, by looking at this diagram, that W must be $\frac{2}{3}$ of a pound. (If you get an epitaph like that on your gravestone, you are doing fine.)

Let us now illustrate the energy principle with a more complicated problem, the screw jack shown in Fig. 4-5. A handle 20 inches long is used to turn the screw, which has 10 threads to the inch. We would like to know how much force would be needed at the handle to lift one ton (2000 pounds). If we want to lift the ton one inch, say, then we must turn the handle around ten times. When it goes around once it goes approximately 126 inches. The handle must thus travel 1260 inches, and if we used various pulleys, etc., we would be lifting our one ton with an unknown smaller weight W applied to the end of the handle. So we find out that W is about 1.6 pounds. This is a result of the conservation of energy.

Take now the somewhat more complicated example shown in Fig. 4-6. A rod or bar, 8 feet long, is supported at one end. In the middle of the bar is a weight of 60 pounds, and at a distance of two feet from the support there is a weight of 100 pounds. How hard do we have to lift the end of the bar in order to keep it balanced, disregarding the weight of the bar? Suppose we put a pulley at one end and hang a weight on the pulley. How big would the weight W have to be in order for it to balance? We imagine that the weight falls any arbitrary distance—to make it easy for ourselves suppose it goes down 4 inches—how high would the two load weights rise? The center rises 2 inches, and the point a quarter of the way from the fixed end lifts 1 inch. Therefore, the principle that the sum of the heights times the weights does not change tells us that the weight W times 4 inches down, plus 60 pounds times 2 inches up, plus 100 pounds times 1 inch has to add up to nothing:

$$-4W + (2)(60) + (1)(100) = 0, \quad W = 55 \text{ lb.} \quad (4.5)$$

Thus we must have a 55-pound weight to balance the bar. In this way we can work out the laws of "balance"—the statics of complicated bridge arrangements, and so on. This approach is called *the principle of virtual work*, because in order to apply this argument we had to *imagine* that the structure moves a little—even though it is not *really* moving or even *movable*. We use the very small imagined motion to apply the principle of conservation of energy.

4.3 Kinetic energy

To illustrate another type of energy we consider a pendulum (Fig. 4-7). If we pull the mass aside and release it, it swings back and forth. In its motion, it loses height in going from either end to the center. Where does the potential energy go? Gravitational energy disappears when it is down at the bottom; nevertheless, it will climb up again. The gravitational energy must have gone into another form. Evidently it is by virtue of its *motion* that it is able to climb up again, so we have the conversion of gravitational energy into some other form when it reaches the bottom.

We must get a formula for the energy of motion. Now, recalling our arguments about reversible machines, we can easily see that in the motion at the bottom must be a quantity of energy which permits it to rise a certain height, and which has nothing to do with the *machinery* by which it comes up or the *path* by which it comes up. So we have an equivalence formula something like the one we wrote for the child's blocks. We have another form to represent the energy. It is easy to say what it is. The kinetic energy at the bottom equals the weight times the height that it could go, corresponding to its velocity: $K.E. = WH$. What we need is the formula which tells us the height by some rule that has to do with the motion of objects. If we start something out with a certain velocity, say straight up, it will reach a certain height; we do not know what it is yet, but it depends on the velocity—there is a formula for that. Then to find the formula for kinetic energy

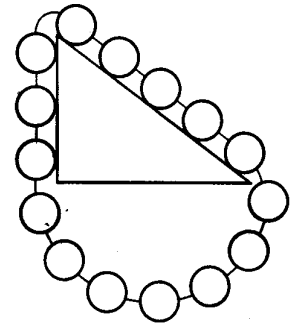


Fig. 4-4. The epitaph of Stevinus.

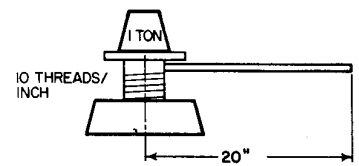


Fig. 4-5. A screw jack.

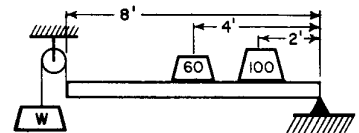


Fig. 4-6. Weighted rod supported on one end.

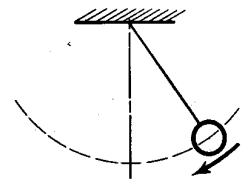


Fig. 4-7. Pendulum.

for an object moving with velocity V , we must calculate the height that it could reach, and multiply by the weight. We shall soon find that we can write it this way:

$$\text{K.E.} = WV^2/2g. \quad (4.6)$$

Of course, the fact that motion has energy has nothing to do with the fact that we are in a gravitational field. It makes no difference *where* the motion came from. This is a general formula for various velocities. Both (4.3) and (4.6) are approximate formulas, the first because it is incorrect when the heights are great, i.e., when the heights are so high that gravity is weakening; the second, because of the relativistic correction at high speeds. However, when we do finally get the exact formula for the energy, then the law of conservation of energy is correct.

4.4 Other forms of energy

We can continue in this way to illustrate the existence of energy in other forms. First, consider elastic energy. If we pull down on a spring, we must do some work, for when we have it down, we can lift weights with it. Therefore in its stretched condition it has a possibility of doing some work. If we were to evaluate the sums of weights times heights, it would not check out—we must add something else to account for the fact that the spring is under tension. Elastic energy is the formula for a spring when it is stretched. How much energy is it? If we let go, the elastic energy, as the spring passes through the equilibrium point, is converted to kinetic energy and it goes back and forth between compressing or stretching the spring and kinetic energy of motion. (There is also some gravitational energy going in and out, but we can do this experiment "sideways" if we like.) It keeps going until the losses—Aha! We have cheated all the way through by putting on little weights to move things or saying that the machines are reversible, or that they go on forever, but we can see that things do stop, eventually. Where is the energy when the spring has finished moving up and down? This brings in *another* form of energy: *heat energy*.

Inside a spring or a lever there are crystals which are made up of lots of atoms, and with great care and delicacy in the arrangement of the parts one can try to adjust things so that as something rolls on something else, none of the atoms do any jiggling at all. But one must be very careful. Ordinarily when things roll, there is bumping and jiggling because of the irregularities of the material, and the atoms start to wiggle inside. So we lose track of that energy; we find the atoms are wiggling inside in a random and confused manner after the motion slows down. There is still kinetic energy, all right, but it is not associated with visible motion. What a dream! How do we *know* there is still kinetic energy? It turns out that with thermometers you can find out that, in fact, the spring or the lever is *warmer*, and that there is really an increase of kinetic energy by a definite amount. We call this form of energy *heat energy*, but we know that it is not really a new form, it is just kinetic energy—internal motion. (One of the difficulties with all these experiments with matter that we do on a large scale is that we cannot really demonstrate the conservation of energy and we cannot really make our reversible machines, because every time we move a large clump of stuff, the atoms do not remain absolutely undisturbed, and so a certain amount of random motion goes into the atomic system. We cannot see it, but we can measure it with thermometers, etc.)

There are many other forms of energy, and of course we cannot describe them in any more detail just now. There is electrical energy, which has to do with pushing and pulling by electric charges. There is radiant energy, the energy of light, which we know is a form of electrical energy because light can be represented as wiggings in the electromagnetic field. There is chemical energy, the energy which is released in chemical reactions. Actually, elastic energy is, to a certain extent, like chemical energy, because chemical energy is the energy of the attraction of the atoms, one for the other, and so is elastic energy. Our modern understanding is the following: chemical energy has two parts, kinetic energy of the electrons inside the atoms, so part of it is kinetic, and electrical energy of interaction of the

electrons and the protons—the rest of it, therefore, is electrical. Next we come to nuclear energy, the energy which is involved with the arrangement of particles inside the nucleus, and we have formulas for that, but we do not have the fundamental laws. We know that it is not electrical, not gravitational, and not purely chemical, but we do not know what it is. It seems to be an additional form of energy. Finally, associated with the relativity theory, there is a modification of the laws of kinetic energy, or whatever you wish to call it, so that kinetic energy is combined with another thing called *mass energy*. An object has energy from its sheer *existence*. If I have a positron and an electron, standing still doing nothing—never mind gravity, never mind anything—and they come together and disappear, radiant energy will be liberated, in a definite amount, and the amount can be calculated. All we need know is the mass of the object. It does not depend on what it is—we make two things disappear, and we get a certain amount of energy. The formula was first found by Einstein; it is $E = mc^2$.

It is obvious from our discussion that the law of conservation of energy is enormously useful in making analyses, as we have illustrated in a few examples without knowing all the formulas. If we had all the formulas for all kinds of energy, we could analyze how many processes should work without having to go into the details. Therefore conservation laws are very interesting. The question naturally arises as to what other conservation laws there are in physics. There are two other conservation laws which are analogous to the conservation of energy. One is called the conservation of linear momentum. The other is called the conservation of angular momentum. We will find out more about these later. In the last analysis, we do not understand the conservation laws deeply. We do not understand the conservation of energy. We do not understand energy as a certain number of little blobs. You may have heard that photons come out in blobs and that the energy of a photon is Planck's constant times the frequency. That is true, but since the frequency of light can be anything, there is no law that says that energy has to be a certain definite amount. Unlike Dennis' blocks, there can be any amount of energy, at least as presently understood. So we do not understand this energy as counting something at the moment, but just as a mathematical quantity, which is an abstract and rather peculiar circumstance. In quantum mechanics it turns out that the conservation of energy is very closely related to another important property of the world, *things do not depend on the absolute time*. We can set up an experiment at a given moment and try it out, and then do the same experiment at a later moment, and it will behave in exactly the same way. Whether this is strictly true or not, we do not know. If we assume that it is true, and add the principles of quantum mechanics, then we can deduce the principle of the conservation of energy. It is a rather subtle and interesting thing, and it is not easy to explain. The other conservation laws are also linked together. The conservation of momentum is associated in quantum mechanics with the proposition that it makes no difference *where* you do the experiment, the results will always be the same. As independence in space has to do with the conservation of momentum, independence of time has to do with the conservation of energy, and finally, if we *turn* our apparatus, this too makes no difference, and so the invariance of the world to angular orientation is related to the conservation of *angular momentum*. Besides these, there are three other conservation laws, that are exact so far as we can tell today, which *are* much simpler to understand because they are in the nature of counting blocks.

The first of the three is the *conservation of charge*, and that merely means that you count how many positive, minus how many negative electrical charges you have, and the number is never changed. You may get rid of a positive with a negative, but you do not create any net excess of positives over negatives. Two other laws are analogous to this one—one is called the *conservation of baryons*. There are a number of strange particles, a neutron and a proton are examples, which are called baryons. In any reaction whatever in nature, if we count how many baryons are coming into a process, the number of baryons* which come out

*Counting antibaryons as — 1 baryon.

will be exactly the same. There is another law, the *conservation of leptons*. We can say that the group of particles called leptons are: electron, mu meson, and neutrino. There is an antielectron which is a positron, that is, a -1 lepton. Counting the total number of leptons in a reaction reveals that the number in and out never changes, at least so far as we know at present.

These are the six conservation laws, three of them subtle, involving space and time, and three of them simple, in the sense of counting something.

With regard to the conservation of energy, we should note that *available* energy is another matter—there is a lot of jiggling around in the atoms of the water of the sea, because the sea has a certain temperature, but it is impossible to get them herded into a definite motion without taking energy from somewhere else. That is, although we know for a fact that energy is conserved, the energy available for human utility is not conserved so easily. The laws which govern how much energy is available are called the *laws of thermodynamics* and involve a concept called entropy for irreversible thermodynamic processes.

Finally, we remark on the question of where we can get our supplies of energy today. Our supplies of energy are from the sun, rain, coal, uranium, and hydrogen. The sun makes the rain, and the coal also, so that all these are from the sun. Although energy is conserved, nature does not seem to be interested in it; she liberates a lot of energy from the sun, but only one part in two billion falls on the earth. Nature has conservation of energy, but does not really care; she spends a lot of it in all directions. We have already obtained energy from uranium; we can also get energy from hydrogen, but at present only in an explosive and dangerous condition. If it can be controlled in thermonuclear reactions, it turns out that the energy that can be obtained from 10 quarts of water per second is equal to all of the electrical power generated in the United States. With 150 gallons of running water a minute, you have enough fuel to supply all the energy which is used in the United States today! Therefore it is up to the physicist to figure out how to liberate us from the need for having energy. It can be done.

Time and Distance

5-1 Motion

In this chapter we shall consider some aspects of the concepts of *time* and *distance*. It has been emphasized earlier that physics, as do all the sciences, depends on *observation*. One might also say that the development of the physical sciences to their present form has depended to a large extent on the emphasis which has been placed on the making of *quantitative* observations. Only with quantitative observations can one arrive at quantitative relationships, which are the heart of physics.

Many people would like to place the beginnings of physics with the work done 350 years ago by Galileo, and to call him the first physicist. Until that time, the study of motion had been a philosophical one based on arguments that could be thought up in one's head. Most of the arguments had been presented by Aristotle and other Greek philosophers, and were taken as "proven." Galileo was skeptical, and did an experiment on motion which was essentially this: He allowed a ball to roll down an inclined trough and observed the motion. He did not, however, just look; he measured *how far* the ball went in *how long a time*.

The way to measure a distance was well known long before Galileo, but there were no accurate ways of measuring time, particularly short times. Although he later devised more satisfactory clocks (though not like the ones we know), Galileo's first experiments on motion were done by using his pulse to count off equal intervals of time. Let us do the same.

We may count off beats of a pulse as the ball rolls down the track: "one . . . two . . . three . . . four . . . five . . . six . . . seven . . . eight..." We ask a friend to make a small mark at the location of the ball at each count; we can then measure the *distance* the ball travelled from the point of release in one, or two, or three, etc., equal intervals of time. Galileo expressed the result of *his* observations in this way: if the location of the ball is marked at 1, 2, 3, 4, . . . units of time from the instant of its release, those marks are distant from the starting point in proportion to the numbers 1, 4, 9, 16, . . . Today we would say the distance is proportional to the square of the time:

$$D \propto t^2.$$

The study of *motion*, which is basic to all of physics, treats with the questions: where? and when?

5-2 Time

Let us consider first what we mean by *time*. What *is* time? It would be nice if we could find a good definition of time. Webster defines "a time" as "a period," and the latter as "a time," which doesn't seem to be very useful. Perhaps we should say: "Time is what happens when nothing else happens." Which also doesn't get us very far. Maybe it is just as well if we face the fact that time is one of the things we probably cannot define (in the dictionary sense), and just say that it is what we already know it to be: it is how long we wait!

What really matters anyway is not how we *define* time, but how we measure it. One way of measuring time is to utilize something which happens over and over again in a regular fashion—something which is *periodic*. For example, a day. A day seems to happen over and over again. But when you begin to think

5-1 Motion

5-2 Time

5-3 Short times

5-4 Long times

5-5 Units and standards of time

5-6 Large distances

5-7 Short distances

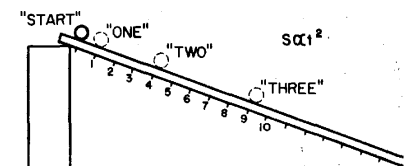


Fig. 5-1. A ball rolls down an inclined track.

about it, you might well ask: "Are days periodic; are they regular? Are all days the same length?" One certainly has the impression that days in summer are longer than days in winter. Of course, some of the days in winter seem to get awfully long if one is very bored. You have certainly heard someone say, "My, but this has been a long day!"

It does seem, however, that days are about the same length *on the average*. Is there any way we can test whether the days are the same length—either from one day to the next, or at least on the average? One way is to make a comparison with some other periodic phenomenon. Let us see how such a comparison might be made with an hour glass. With an hour glass, we can "create" a periodic occurrence if we have someone standing by it day and night to turn it over whenever the last grain of sand runs out.

We could then count the turnings of the glass from each morning to the next. We would find, this time, that the number of "hours" (i.e., turnings of the glass) was not the same each "day." We should distrust the sun, or the glass, or both. After some thought, it might occur to us to count the "hours" from noon to noon. (Noon is here defined *not* as 12:00 o'clock, but that instant when the sun is at its highest point.) We would find, this time, that the number of "hours" each day is the same.

We now have some confidence that both the "hour" and the "day" have a regular periodicity, i.e., mark off successive equal intervals of time, although we have not *proved* that either one is "really" periodic. Someone might question whether there might not be some omnipotent being who would slow down the flow of sand every night and speed it up during the day. Our experiment does not, of course, give us an answer to this sort of question. All we can say is that we find that a regularity of one kind fits together with a regularity of another kind. We can just say that we base our *definition* of time on the repetition of some apparently periodic event.

5-3 Short times

We should now notice that in the process of checking on the reproducibility of the day, we have received an important by-product. We have found a way of measuring, more accurately, *fractions* of a day. We have found a way of counting time in smaller pieces. Can we carry the process further, and learn to measure even smaller intervals of time?

Galileo decided that a given pendulum always swings back and forth in equal intervals of time so long as the size of the swing is kept small. A test comparing the number of swings of a pendulum in one "hour" shows that such is indeed the case. We can in this way mark fractions of an hour. If we use a mechanical device to count the swings—and to keep them going—we have the pendulum clock of our grandfathers.

Let us agree that if our pendulum oscillates 3600 times in one hour (and if there are 24 such hours in a day), we shall call each period of the pendulum one "second." We have then divided our original unit of time into approximately 10⁵ parts. We can apply the same principles to divide the second into smaller and smaller intervals. It is, you will realize, not practical to make mechanical pendulums which go arbitrarily fast, but we can now make *electrical* pendulums, called oscillators, which can provide a periodic occurrence with a very short period of swing. In these electronic oscillators it is an electrical current which swings to and fro, in a manner analogous to the swinging of the bob of the pendulum.

We can make a series of such electronic oscillators, each with a period 10 times shorter than the previous one. We may "calibrate" each oscillator against the next slower one by counting the number of swings it makes for one swing of the slower oscillator. When the period of oscillation of our clock is shorter than a fraction of a second, we cannot count the oscillations without the help of some device which extends our powers of observation. One such device is the electron-beam oscilloscope, which acts as a sort of microscope for short times. This device plots on a fluorescent screen a graph of electrical current (or voltage) versus time.

By connecting the oscilloscope to two of our oscillators in sequence, so that it plots a graph first of the current in one of our oscillators and then of the current in the other, we get two graphs like those shown in Fig. 5-2. We can readily determine the number of periods of the faster oscillator in one period of the slower oscillator.

With modern electronic techniques, oscillators have been built with periods as short as about 10^{-12} second, and they have been calibrated (by comparison methods such as we have described) in terms of our standard unit of time, the second. With the invention and perfection of the "laser," or light amplifier, in the past few years, it has become possible to make oscillators with even shorter periods than 10^{-12} second, but it has not yet been possible to calibrate them by the methods which have been described, although it will no doubt soon be possible.

Times shorter than 10^{-12} second have been measured, but by a different technique. In effect, a different *definition* of "time" has been used. One way has been to observe the *distance* between two happenings on a moving object. If, for example, the headlights of a moving automobile are turned on and then off, we can figure out *how long* the lights were on if we know *where* they were turned on and off and how fast the car was moving. The time is the distance over which the lights were on divided by the speed.

Within the past few years, just such a technique was used to measure the lifetime of the π^0 -meson. By observing in a microscope the minute tracks left in a photographic emulsion in which π^0 -mesons had been created one saw that a π^0 -meson (known to be travelling at a certain speed nearly that of light) went a distance of about 10^{-7} meter, on the average, before disintegrating. It lived for only about 10^{-16} sec. It should be emphasized that we have here used a somewhat different definition of "time" than before. So long as there are no inconsistencies in our understanding, however, we feel fairly confident that our definitions are sufficiently equivalent.

By extending our techniques—and if necessary our definitions—still further we can infer the time duration of still faster physical events. We can speak of the period of a nuclear vibration. We can speak of the lifetime of the newly discovered strange resonances (particles) mentioned in Chapter 2. Their complete life occupies a time span of only 10^{-24} second, approximately the time it would take light (which moves at the fastest known speed) to cross the nucleus of hydrogen (the smallest known object).

What about still smaller times? Does "time" exist on a still smaller scale? Does it make any sense to speak of smaller times if we cannot measure—or perhaps even think sensibly about—something which happens in a shorter time? Perhaps not. These are some of the open questions which you will be asking and perhaps answering in the next twenty or thirty years.

5-4 Long times

Let us now consider times longer than one day. Measurement of longer times is easy; we just count the days—so long as there is someone around to do the counting. First we find that there is another natural periodicity: the year, about 365 days. We have also discovered that nature has sometimes provided a counter for the years, in the form of tree rings or river-bottom sediments. In some cases we can use these natural time markers to determine the time which has passed since some early event.

When we cannot count the years for the measurement of long times, we must look for other ways to measure. One of the most successful is the use of radioactive material as a "clock." In this case we do not have a periodic occurrence, as for the day or the pendulum, but a new kind of "regularity." We find that the radioactivity of a particular sample of material decreases by the same *fraction* for successive equal increases in its age. If we plot a graph of the radioactivity observed as a function of time (say in days), we obtain a curve like that shown in Fig. 5-3. We observe that if the radioactivity decreases to one-half in T days (called the "half-life"), then it decreases to one-quarter in another T days, and so

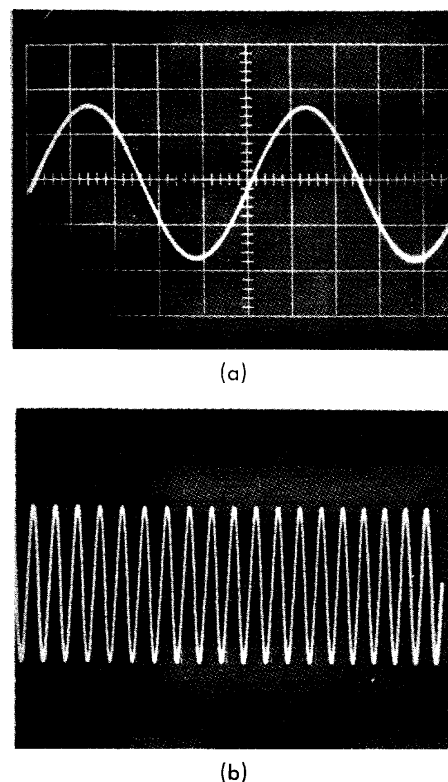


Fig. 5-2. Two views of an oscilloscope screen. In (a) the oscilloscope is connected to one oscillator, in (b) it is connected to an oscillator with a period one-tenth as long.

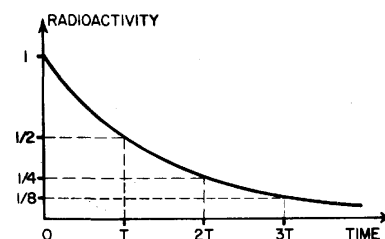


Fig. 5-3. The decrease with time of radioactivity. The activity decreases by one-half in each "half-life," T .

		TIMES								MEAN LIFE OF
YEARS	SECONDS									
		? ? ? ? ? ? ? ?								
	10^{18}	Age of universe								
10^9		Age of earth								U^{238}
	10^{15}									
10^6		Earliest men								
	10^{12}	Age of pyramids								Ra^{226}
10^3		Age of U.S.								
	10^9	Life of a man								H^3
1										
	10^6	One day								
	10^3	Light goes from sun to earth								Neutron
	1	One heart beat								
	10^{-3}	Period of sound wave								
	10^{-6}	Period of radiowave								Muon π^\pm -meson
	10^{-9}	Light travels one foot								
	10^{-12}	Period of molecular rotation								
	10^{-15}	Period of atomic vibration								π^0 -meson
	10^{-18}	Light crosses an atom								
	10^{-21}									
		Period of nuclear vibration								
	10^{-24}	Light crosses a nucleus								Strange particle
		? ? ? ? ? ? ? ?								

on. In an arbitrary time interval t there are t/T "half-lives," and the fraction left after this time t is $(\frac{1}{2})^{t/T}$.

If we knew that a piece of material, say a piece of wood, had contained an amount A of radioactive material when it was formed, and we found out by a direct measurement that it now contains the amount B , we could compute the age of the object, t , by solving the equation

$$(\frac{1}{2})^{t/T} = B/A.$$

There are, fortunately, cases in which we can know the amount of radioactivity that was in an object when it was formed. We know, for example, that the carbon dioxide in the air contains a certain small fraction of the radioactive carbon isotope C^{14} (replenished continuously by the action of cosmic rays). If we measure the *total* carbon content of an object, we know that a certain fraction of that amount was originally the radioactive C^{14} ; we know, therefore, the starting amount A to use in the formula above. Carbon-14 has a half-life of 5000 years. By careful measurements we can measure the amount left after 20 half-lives or so and can therefore "date" organic objects which grew as long as 100,000 years ago.

We would like to know, and we think we do know, the life of still older things. Much of our knowledge is based on the measurements of other radioactive isotopes which have different half-lives. If we make measurements with an isotope with a longer half-life, then we are able to measure longer times. Uranium, for example, has an isotope whose half-life is about 10^9 years, so that if some material was formed with uranium in it 10^9 years ago, only half the uranium would remain today. When the uranium disintegrates, it changes into lead. Consider a piece of rock which was formed a long time ago in some chemical process. Lead, being of a chemical nature different from uranium, would appear in one part of the rock and uranium would appear in another part of the rock. The uranium and lead

would be separate. If we look at that piece of rock today, where there should only be uranium we will now find a certain fraction of uranium and a certain fraction of lead. By comparing these fractions, we can tell what percent of the uranium disappeared and changed into lead. By this method, the age of certain rocks has been determined to be several billion years. An extension of this method, not using particular rocks but looking at the uranium and lead in the oceans and using averages over the earth, has been used to determine (within the past few years) that the age of the earth itself is approximately 5.5 billion years.

It is encouraging that the age of the earth is found to be the same as the age of the meteorites which land on the earth, as determined by the uranium method. It appears that the earth was formed out of rocks floating in space, and that the meteorites are, quite likely, some of that material left over. At some time more than five billion years ago, the universe started. It is now believed that at least our part of the universe had its beginning about ten or twelve billion years ago. We do not know what happened before then. In fact, we may well ask again: Does the question make any sense? Does an earlier time have any meaning?

5-5 Units and standards of time

We have implied that it is convenient if we start with some standard unit of time, say a day or a second, and refer all other times to some multiple or fraction of this unit. What shall we take as our basic standard of time? Shall we take the human pulse? If we compare pulses, we find that they seem to vary a lot. On comparing two clocks, one finds they do not vary so much. You might then say, well, let us take a clock. But whose clock? There is a story of a Swiss boy who wanted all of the clocks in his town to ring noon at the same time. So he went around trying to convince everyone of the value of this. Everyone thought it was a marvelous idea so long as all of the other clocks rang noon when his did! It is rather difficult to decide whose clock we should take as a standard. Fortunately, we all share one clock—the earth. For a long time the rotational period of the earth has been taken as the basic standard of time. As measurements have been made more and more precise, however, it has been found that the rotation of the earth is not exactly periodic, when measured in terms of the best clocks. These "best" clocks are those which we have reason to believe are accurate because they agree with each other. We now believe that, for various reasons, some days are longer than others, some days are shorter, and on the average the period of the earth becomes a little longer as the centuries pass.

Until very recently we had found nothing much better than the earth's period, so all clocks have been related to the length of the day, and the second has been defined as $1/86400$ of an average day. Recently we have been gaining experience with some natural oscillators which we now believe would provide a more constant time reference than the earth, and which are also based on a natural phenomenon available to everyone. These are the so-called "atomic clocks." Their basic internal period is that of an atomic vibration which is very insensitive to the temperature or any other external effects. These clocks keep time to an accuracy of one part in 10^9 or better. Within the past two years an improved atomic clock which operates on the vibration of the hydrogen atom has been designed and built by Professor Norman Ramsey at Harvard University. He believes that this clock might be 100 times more accurate still. Measurements now in progress will show whether this is true or not.

We may expect that since it has been possible to build clocks much more accurate than astronomical time, there will soon be an agreement among scientists to define the unit of time in terms of one of the atomic clock standards.

5-6 Large distances

Let us now turn to the question of *distance*. How far, or how big, are things? Everybody knows that the way you measure distance is to start with a stick and count. Or start with a thumb and count. You begin with a unit and count. How

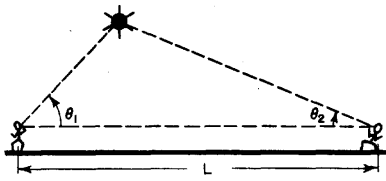


Fig. 5-4. The height of a Sputnik is determined by triangulation

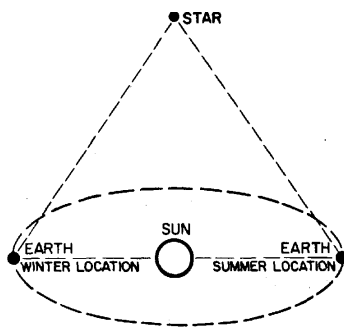


Fig. 5-5. The distance of nearby stars can be measured by triangulation, using the diameter of the earth's orbit as a baseline.

does one measure smaller things? How does one subdivide distance? In the same way that we subdivided time: we take a smaller unit and count the number of such units it takes to make up the longer unit. So we can measure smaller and smaller lengths.

But we do not always mean by distance what one gets by counting off with a meter stick. It would be difficult to measure the horizontal distance between two mountain tops using only a meter stick. We have found by experience that distance can be measured in another fashion: by triangulation. Although this means that we are really using a different definition of distance, when they can both be used they agree with each other. Space is more or less what Euclid thought it was, so the two types of definitions of distance agree. Since they do agree on the earth it gives us some confidence in using triangulation for still larger distances. For example, we were able to use triangulation to measure the height of the first Sputnik. We found that it was roughly 5×10^5 meters high. By more careful measurements the distance to the moon can be measured in the same way. Two telescopes at different places on the earth can give us the two angles we need. It has been found in this way that the moon is 4×10^8 meters away.

We cannot do the same with the sun, or at least no one has been able to yet. The accuracy with which one can focus on a given point on the sun and with which one can measure angles is not good enough to permit us to measure the distance to the sun. Then how can we measure the distance to the sun? We must invent an extension of the idea of triangulation. We measure the relative distances of all the planets by astronomical observations of where the planets appear to be, and we get a picture of the solar system with the proper *relative* distances of everything, but with no *absolute* distance. One absolute measurement is then required, which has been obtained in a number of ways. One of the ways, which was believed until recently to be the most accurate, was to measure the distance from the earth to Eros, one of the small planetoids which passes near the earth every now and then. By triangulation on this little object, one could get the one required scale measurement. Knowing the relative distances of the rest, we can then tell the distance, for example, from the earth to the sun, or from the earth to Pluto.

Within the past year there has been a big improvement in our knowledge of the scale of the solar system. At the Jet Propulsion Laboratory the distance from the earth to Venus was measured quite accurately by a direct radar observation. This, of course, is a still different type of inferred distance. We say we know the speed at which light travels (and therefore, at which radar waves travel), and we assume that it is the same speed everywhere between the earth and Venus. We send the radio wave out, and count the time until the reflected wave comes back. From the *time* we infer a *distance*, assuming we know the speed. We have really another definition of a measurement of distance.

How do we measure the distance to a star, which is much farther away? Fortunately, we can go back to our triangulation method, because the earth moving around the sun gives us a large baseline for measurements of objects outside the solar system. If we focus a telescope on a star in summer and in winter, we might hope to determine these two angles accurately enough to be able to measure the distance to a star.

What if the stars are too far away for us to use triangulation? Astronomers are always inventing new ways of measuring distance. They find, for example, that they can estimate the size and brightness of a star by its color. The color and brightness of many nearby stars—whose distances are known by triangulation—have been measured, and it is found that there is a smooth relationship between the color and the intrinsic brightness of stars (in most cases). If one now measures the color of a distant star, one may use the color-brightness relationship to determine the intrinsic brightness of the star. By measuring how bright the star *appears* to us at the earth (or perhaps we should say how *dim* it appears), we can compute how far away it is. (For a given intrinsic brightness, the apparent brightness decreases with the square of the distance.) A nice confirmation of the correctness of this method of measuring stellar distances is given by the results obtained for groups of stars known as globular clusters. A photograph of such a group is

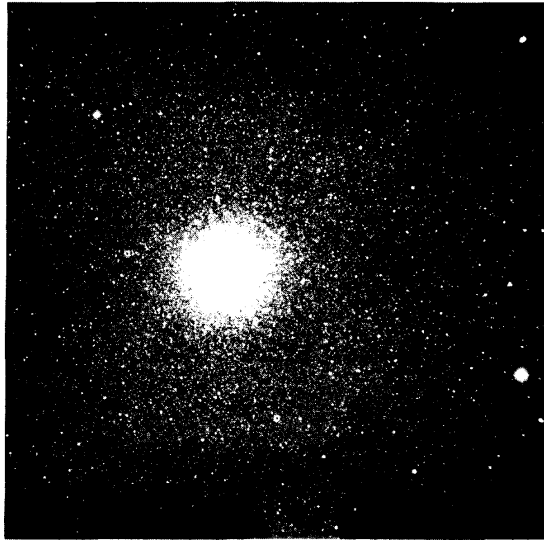


Fig. 5-6. A cluster of stars near the center of our galaxy. Their distance from the earth is 30,000 light-years, or about 3×10^{20} meters.

shown in Fig. 5-6. Just from looking at the photograph one is convinced that these stars are all together. The same result is obtained from distance measurements by the color-brightness method.

A study of many globular clusters gives another important bit of information. It is found that there is a high concentration of such clusters in a certain part of the sky and that most of them are about the same distance from us. Coupling this information with other evidence, we conclude that this concentration of clusters marks the center of our galaxy. We then know the distance to the center of the galaxy—about 10^{20} meters.

Knowing the size of our own galaxy, we have a key to the measurement of still larger distances—the distances to other galaxies. Figure 5-7 is a photograph of a galaxy, which has much the same shape as our own. Probably it is the same size, too. (Other evidence supports the idea that galaxies are all about the same size.) If it is the same size as ours, we can tell its distance. We measure the angle it subtends in the sky; we know its diameter, and we compute its distance—triangulation again!



Fig. 5-7. A spiral galaxy like our own. Presuming that its diameter is similar to that of our own galaxy, we may compute its distance from its apparent size. It is 30 million light-years (3×10^{23} meters) from the earth.

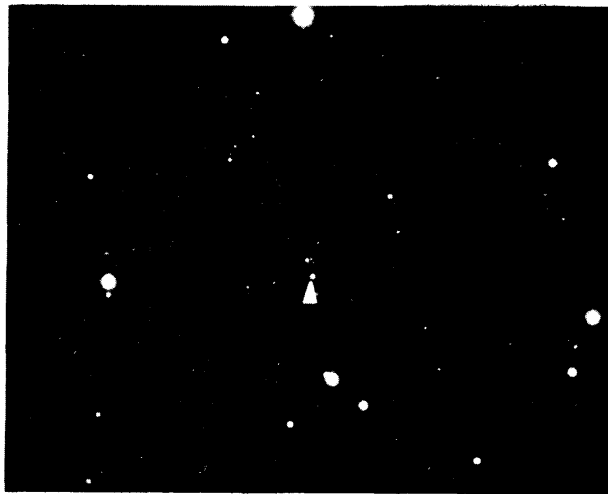


Fig. 5-8. The most distant object, 3C295 in BOOTES (indicated by the arrow), measured by the 200-inch telescope to date (1960).

Photographs of exceedingly distant galaxies have recently been obtained with the giant Palomar telescope. One is shown in Fig. 5-8. It is now believed that some of these galaxies are about halfway to the limit of the universe— 10^{26} meters away—the largest distance we can contemplate!

5-7 Short distances

Now let's think about smaller distances. Subdividing the meter is easy. Without much difficulty we can mark off one thousand equal spaces which add up to one meter. With somewhat more difficulty, but in a similar way (using a good microscope), we can mark off a thousand equal subdivisions of the millimeter to make a scale of microns (millionths of a meter). It is difficult to continue to smaller scales, because we cannot "see" objects smaller than the wavelength of visible light (about 5×10^{-7} meter).

We need not stop, however, at what we can see. With an electron microscope, we can continue the process by making photographs on a still smaller scale, say down to 10^{-8} meter (Fig. 5-9). By indirect measurements—by a kind of triangulation on a microscopic scale—we can continue to measure to smaller and smaller scales. First, from an observation of the way light of short wavelength (x-radiation) is reflected from a pattern of marks of known separation, we determine the wave-

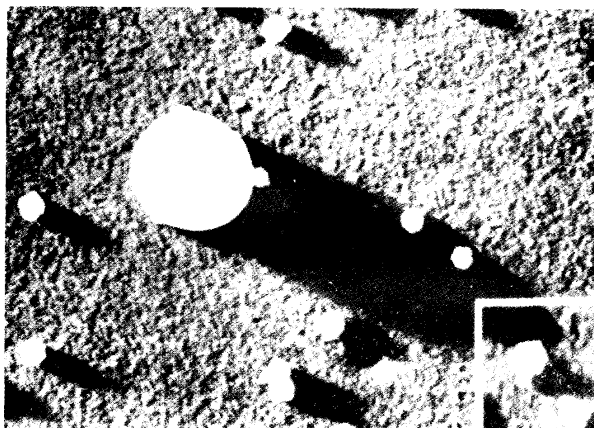


Fig. 5-9. Electron micrograph of some virus molecules. The "large" sphere is for calibration and is known to have a diameter of 2×10^{-7} meter (2000 Å).

DISTANCES		
LIGHT-YEARS	METERS	
	10^{27}	?? ? ? ? ? ? ?
		Edge of universe
10^9	10^{24}	
10^6	10^{21}	To nearest neighbor galaxy
		To center of our galaxy
10^3	10^{18}	
		To nearest star
1	10^{15}	
		Radius of orbit of Pluto
	10^{12}	
		To the sun
	10^9	
		To the moon
	10^6	
		Height of a Sputnik
	10^3	
		Height of a TV antenna tower
	1	Height of a child
	10^{-3}	
		A grain of salt
	10^{-6}	
		A virus
	10^{-9}	
		Radius of an atom
	10^{-12}	
	10^{-15}	Radius of a nucleus
		?? ? ? ? ? ? ?

length of the light vibrations. Then, from the pattern of the scattering of the same light from a crystal, we can determine the relative location of the atoms in the crystal, obtaining results which agree with the atomic spacings also determined by chemical means. We find in this way that atoms have a diameter of about 10^{-10} meter.

There is a large "gap" in physical sizes between the typical atomic dimension of about 10^{-10} meter and the nuclear dimensions 10^{-15} meter, 10^{-5} times smaller. For nuclear sizes, a different way of measuring size becomes convenient. We measure the *apparent area*, σ , called the effective *cross section*. If we wish the radius, we can obtain it from $\sigma = \pi r^2$, since nuclei are nearly spherical.

Measurement of a nuclear cross section can be made by passing a beam of high-energy particles through a thin slab of material and observing the number of particles which do not get through. These high-energy particles will plow right through the thin cloud of electrons and will be stopped or deflected only if they hit the concentrated weight of a nucleus. Suppose we have a piece of material 1 centimeter thick. There will be about 10^8 atomic layers. But the nuclei are so small that there is little chance that any nucleus will lie behind another. We might *imagine* that a highly magnified view of the situation—looking along the particle beam—would look like Fig. 5-10.

The chance that a very small particle will hit a nucleus on the trip through is just the total area covered by the profiles of the nuclei divided by the total area in the picture. Suppose that we know that in an area A of our slab of material there are N atoms (each with one nucleus, of course). Then the total area "covered" by the nuclei is $N\sigma/A$. Now let the number of particles of our beam which arrive at the slab be n_1 and the number which come out the other side be n_2 . The fraction which do *not* get through is $(n_1 - n_2)/n_1$, which should just equal the

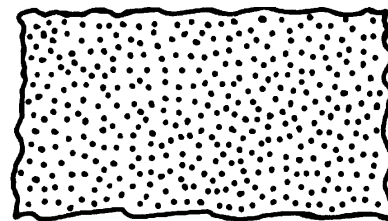


Fig. 5-10. Imagined view through a block of carbon 1 cm thick if only the nuclei were observed.

fraction of the area covered. We can obtain the radius of the nucleus from the equation*

$$\pi r^2 = \sigma = \frac{A}{N} \frac{n_1 - n_2}{n_1}.$$

From such an experiment we find that the radii of the nuclei are from about 1 to 6 times 10^{-15} meter. The length unit 10^{-15} meter is called *the fermi*, in honor of Enrico Fermi (1901-1958).

What do we find if we go to smaller distances? Can we measure smaller distances? Such questions are not yet answerable. It has been suggested that the still unsolved mystery of nuclear forces may be unravelled only by some modification of our idea of space, or measurement, at such small distances.

It might be thought that it would be a good idea to use some natural length as our unit of length—say the radius of the earth or some fraction of it. The meter was originally intended to be such a unit and was defined to be $(\pi/2) \times 10^{-7}$ times the earth's radius. It is neither convenient nor very accurate to determine the unit of length in this way. For a long time it has been agreed internationally that the meter would be defined as the distance between two scratches on a bar kept in a special laboratory in France. More recently, it has been realized that this definition is neither as precise as would be useful, nor as permanent or universal as one would like. It is currently being considered that a new definition be adopted, an agreed-upon (arbitrary) number of wavelengths of a chosen spectral line.

Measurements of distance and of time give results which depend on the observer. Two observers moving with respect to each other will not measure the same distances and times when measuring what appear to be the same things. Distances and time intervals have different magnitudes, depending on the coordinate system (or "frame of reference") used for making the measurements. We shall study this subject in more detail in a later chapter.

Perfectly precise measurements of distances or times are not permitted by the laws of nature. We have mentioned earlier that the errors in a measurement of the position of an object must be at least as large as

$$\Delta x = h/\Delta p,$$

where h is a small quantity called "Planck's constant" and Δp is the error in our knowledge of the momentum (mass times velocity) of the object whose position we are measuring. It was also mentioned that the uncertainty in position measurements is related to the wave nature of particles.

The relativity of space and time implies that time measurements have also a minimum error, given in fact by

$$\Delta t = h/\Delta E,$$

where ΔE is the error in our knowledge of the energy of the process whose time period we are measuring. If we wish to know *more* precisely *when* something happened we must know less about *what* happened, because our knowledge of the energy involved will be less. The time uncertainty is also related to the wave nature of matter.

* This equation is right only if the area covered by the nuclei is a small fraction of the total, i.e., if $(n_1 - n_2)/n_1$ is much less than 1. Otherwise we must make a correction for the fact that some nuclei will be partly obscured by the nuclei in front of them.

Probability

"The true logic of this world is in the calculus of probabilities."

—James Clerk Maxwell

6-1 Chance and likelihood

"Chance" is a word which is in common use in everyday living. The radio reports speaking of tomorrow's weather may say: "There is a sixty percent chance of rain." You might say: "There is a small chance that I shall live to be one hundred years old." Scientists also use the word chance. A seismologist may be interested in the question: "What is the chance that there will be an earthquake of a certain size in Southern California next year?" A physicist might ask the question: "What is the chance that a particular geiger counter will register twenty counts in the next ten seconds?" A politician or statesman might be interested in the question: "What is the chance that there will be a nuclear war within the next ten years?" You may be interested in the chance that you will learn something from this chapter.

By *chance*, we mean something like a guess. Why do we make guesses? We make guesses when we wish to make a judgment but have incomplete information or uncertain knowledge. We want to make a guess as to what things are, or what things are likely to happen. Often we wish to make a guess because we have to make a decision. For example: Shall I take my raincoat with me tomorrow? For what earth movement should I design a new building? Shall I build myself a fallout shelter? Shall I change my stand in international negotiations? Shall I go to class today?

Sometimes we make guesses because we wish, with our limited knowledge, to say as much as we *can* about some situation. Really, any generalization is in the nature of a guess. Any physical theory is a kind of guesswork. There are good guesses and there are bad guesses. The theory of probability is a system for making better guesses. The language of probability allows us to speak quantitatively about some situation which may be highly variable, but which does have some consistent average behavior.

Let us consider the flipping of a coin. If the toss—and the coin—are "honest," we have no way of knowing what to expect for the outcome of any particular toss. Yet we would feel that in a large number of tosses there should be about equal numbers of heads and tails. We say: "The probability that a toss will land heads is 0.5."

We speak of probability only for observations that we contemplate being made in the future. By the "*probability*" of a particular outcome of an observation we mean our estimate for the most likely fraction of a number of repeated observations that will yield that particular outcome. If we imagine repeating an observation—such as looking at a freshly tossed coin— N times, and if we call N_A our estimate of the most likely number of our observations that will give some specified result A , say the result "heads," then by $P(A)$, the probability of observing A , we mean

$$P(A) = N_A/N. \quad (6.1)$$

Our definition requires several comments. First of all, we may speak of a probability of something happening only if the occurrence is a possible outcome of some *repeatable* observation. It is not clear that it would make any sense to ask: "What is the probability that there is a ghost in that house?"

6-1 Chance and likelihood

6-2 Fluctuations

6-3 The random walk

6-4 A probability distribution

6-5 The uncertainty principle

You may object that no situation is *exactly* repeatable. That is right. Every different observation must at least be at a different time or place. All we can say is that the “repeated” observations should, for our intended purposes, *appear to be equivalent*. We should assume, at least, that each observation was made from an equivalently prepared situation, and especially with the same degree of ignorance at the start. (If we sneak a look at an opponent’s hand in a card game, our estimate of our chances of winning are different than if we do not!)

We should emphasize that N and N_A in Eq. (6.1) are *not* intended to represent numbers based on actual observations. N_A is our best *estimate* of what *would* occur in N *imagined* observations. Probability depends, therefore, on our knowledge and on our ability to make estimates. In effect, on our common sense! Fortunately, there is a certain amount of agreement in the common sense of many things, so that different people will make the same estimate. Probabilities need not, however, be “absolute” numbers. Since they depend on our ignorance, they may become different if our knowledge changes.

You may have noticed another rather “subjective” aspect of our definition of probability. We have referred to N_A as “our estimate of the most likely number . . .” We do not mean that we expect to observe *exactly* N_A , but that we expect a number *near* N_A , and that the number N_A is *more likely* than any other number in the vicinity. If we toss a coin, say, 30 times, we should expect that the number of heads would not be very likely to be exactly 15, but rather only some number near to 15, say 12, 13, 14, 15, 16, or 17. However, if we *must* choose, we would decide that 15 heads is *more likely* than any other number. We would write $P(\text{heads}) = 0.5$.

Why did we choose 15 as more likely than any other number? We must have argued with ourselves in the following manner: If the most likely number of heads is N_H in a total number of tosses N , then the most likely number of tails N_T is $(N - N_H)$. (We are assuming that every toss gives *either* heads *or* tails, and no “other” result!) But if the coin is “honest,” there is no preference for heads or tails. Until we have some reason to think the coin (or toss) is dishonest, we must give equal likelihoods for heads and tails. So we must set $N_T = N_H$. It follows that $N_T = N_H = N/2$, or $P(H) = p(T) = 0.5$.

We can generalize our reasoning to *any* situation in which there are m different but “equivalent” (that is, equally likely) possible results of an observation. If an observation can yield m different results, and we have reason to believe that any one of them is as likely as any other, then the probability of a *particular* outcome A is $P(A) = 1/m$.

If there are seven different-colored balls in an opaque box and we pick one out “at random” (that is, without looking), the probability of getting a ball of a particular color is $\frac{1}{7}$. The probability that a “blind draw” from a shuffled deck of 52 cards will show the ten of hearts is $\frac{1}{52}$. The probability of throwing a double-one with dice is $\frac{1}{36}$.

In Chapter 5 we described the size of a nucleus in terms of its apparent area, or “cross section.” When we did so we were really talking about probabilities. When we shoot a high-energy particle at a thin slab of material, there is some chance that it will pass right through and some chance that it will hit a nucleus. (Since the nucleus is so small that we cannot *see* it, we cannot aim right at a nucleus. We must “shoot blind.”) If there are n atoms in our slab and the nucleus of each atom has a cross-sectional area σ , then the total area “shadowed” by the nuclei is $n\sigma$. In a large number N of random shots, we expect that the number of hits N_C of *some* nucleus will be in the ratio to N as the shadowed area is to the total area of the slab:

$$N_C/N = n\sigma/A. \quad (6.2)$$

We may say, therefore, that the *probability* that any one projectile particle will suffer a collision in passing through the slab is

$$P_C = \frac{n}{A} \sigma, \quad (6.3)$$

where n/A is the number of atoms per unit area in our slab.

6-2 Fluctuations

We would like now to use our ideas about probability to consider in some greater detail the question: "How many heads do I really *expect* to get if I toss a coin N times?" Before answering the question, however, let us look at what does happen in such an "experiment." Figure 6-1 shows the results obtained in the first three "runs" of such an experiment in which $N = 30$. The sequences of "heads" and "tails" are shown just as they were obtained. The first game gave 11 heads; the second also 11; the third 16. In three trials we did not once get 15 heads. Should we begin to suspect the coin? Or were we wrong in thinking that the most likely number of "heads" in such a game is 15? Ninety-seven more runs were made to obtain a total of 100 experiments of 30 tosses each. The results of the experiment are given in Table 6-1.*

```

H   x x x x       x x x x x x x x   11
T   x x x x x x x x x x x x x x   19
H   x   x   x x x x x   x x x x   11
T   x x x x x x x x x   x x x x x x   19
H   x   x x x x x x   x x x x x x x   16
T   x x x   x x x x x   x x x x x x   14

```

Fig. 6-1. Observed sequences of heads and tails in three games of 30 tosses each.

Table 6-1

Number of heads in successive trials of 30 tosses of a coin.

11	16	17	15	17	16	19	18	15	13
11	17	17	12	20	23	11	16	17	14
16	12	15	10	18	17	13	15	14	15
16	12	11	22	12	20	12	15	16	12
16	10	15	13	14	16	15	16	13	18
14	14	13	16	15	19	21	14	12	15
16	11	16	14	17	14	11	16	17	16
19	15	14	12	18	15	14	21	11	16
17	17	12	13	14	17	9	13	16	13

90 trials

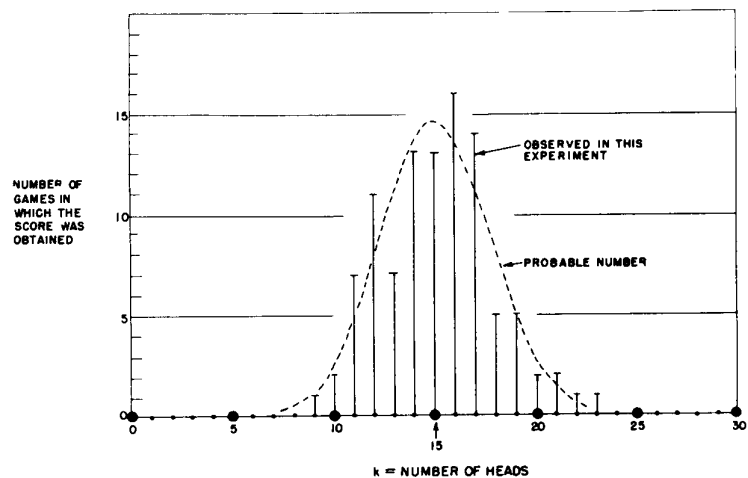


Fig. 6-2. Summary of the results of 100 games of 30 tosses each. The vertical bars show the number of games in which a score of k heads was obtained. The dashed curve shows the expected numbers of games with the score k obtained by a probability computation.

Looking at the numbers in Table 6-1, we see that most of the results are "near" 15, in that they are between 12 and 18. We can get a better feeling for the details of these results if we plot a graph of the *distribution* of the results. We count the number of games in which a score of k was obtained, and plot this number for each k . Such a graph is shown in Fig. 6-2. A score of 15 heads was obtained in 13 games. A score of 14 heads was also obtained 13 times. Scores of 16 and 17 were each obtained *more* than 13 times. Are we to conclude that there is some bias toward heads? Was our "best estimate" not good enough? Should

* After the first three games, the experiment was actually done by shaking 30 pennies violently in a box and then counting the number of heads that showed.

we conclude now that the “most likely” score for a run of 30 tosses is really 16 heads? But wait! In all the games taken together, there were 3000 tosses. And the total number of heads obtained was 1492. The fraction of tosses that gave heads is 0.497, very nearly, but slightly *less* than half. We should certainly *not* assume that the probability of throwing heads is greater than 0.5! The fact that one *particular* set of observations gave 16 heads most often, is a *fluctuation*. We still expect that the *most likely* number of heads is 15.

We may ask the question: “What is the probability that a game of 30 tosses will yield 15 heads—or 16, or any other number?” We have said that in a game of one toss, the probability of obtaining *one* head is 0.5, and the probability of obtaining no head is 0.5. In a game of two tosses there are *four* possible outcomes: *HH, HT, TH, TT*. Since each of these sequences is equally likely, we conclude that (a) the probability of a score of two heads is $\frac{1}{4}$, (b) the probability of a score of one head is $\frac{2}{4}$, (c) the probability of a zero score is $\frac{1}{4}$. There are *two* ways of obtaining one head, but only one of obtaining either zero or two heads.

Consider now a game of 3 tosses. The third toss is equally likely to be heads or tails. There is only one way to obtain 3 heads: we *must* have obtained 2 heads on the first two tosses, and then heads on the last. There are, however, *three* ways of obtaining 2 heads. We could throw tails after having thrown two heads (one way) or we could throw heads after throwing only one head in the first two tosses (two ways). So for scores of 3-H, 2-H, 1-H, 0-H we have that the number of equally likely ways is 1, 3, 3, 1, with a total of 8 different possible sequences. The probabilities are $\frac{1}{8}, \frac{3}{8}, \frac{3}{8}, \frac{1}{8}$.

The argument we have been making can be summarized by a diagram like that in Fig. 6-3. It is clear how the diagram should be continued for games with a larger number of tosses. Figure 6-4 shows such a diagram for a game of 6 tosses. The number of “ways” to any point on the diagram is just the number of different “paths” (sequences of heads and tails) which can be taken from the starting point. The vertical position gives us the total number of heads thrown. The set of numbers which appears in such a diagram is known as *Pascal’s triangle*. The numbers are also known as the *binomial coefficients*, because they also appear in the expansion of $(a + b)^n$. If we call n the number of tosses and k the number of heads thrown, then the numbers in the diagram are usually designated by the symbol $\binom{n}{k}$. We may remark in passing that the binomial coefficients can also be computed from

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}, \quad (6.4)$$

where $n!$, called “ n -factorial,” represents the product $(n)(n-1)(n-2) \dots (3)(2)(1)$.

We are now ready to compute the probability $P(k, n)$ of throwing k heads in n tosses, using our definition Eq. (6.1). The total number of possible sequences is 2^n (since there are 2 outcomes for each toss), and the number of ways of obtaining k heads is $\binom{n}{k}$, all equally likely, so we have

$$P(k, n) = \frac{\binom{n}{k}}{2^n}. \quad (6.5)$$

Since $P(k, n)$ is the fraction of games which we expect to yield k heads, then in 100 games we should expect to find k heads $100 \cdot P(k, n)$ times. The dashed curve in Fig. 6-2 passes through the points computed from $100 \cdot P(k, 30)$. We see that we *expect* to obtain a score of 15 heads in 14 or 15 games, whereas this score was observed in 13 games. We *expect* a score of 16 in 13 or 14 games, but we obtained that score in 16 games. Such fluctuations are “part of the game.”

The method we have just used can be applied to the most general situation in which there are only two possible outcomes of a single observation. Let us designate the two outcomes by W (for “win”) and L (for “lose”). In the general case, the probability of W or L in a single event need not be equal. Let p be the probability of obtaining the result W . Then q , the probability of L , is necessarily

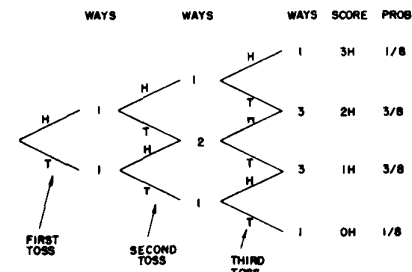


Fig. 6-3. A diagram for showing the number of ways a score of 0, 1, 2, or 3 heads can be obtained in a game of 3 tosses.

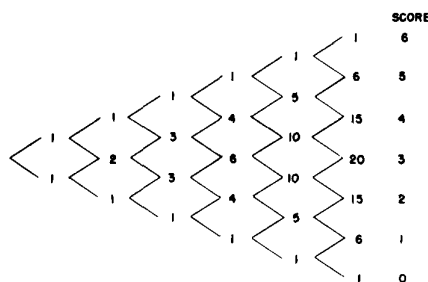


Fig. 6-4. A diagram like that of Fig. 6-3, for a game of 6 tosses.

$(1 - p)$. In a set of n trials, the probability $P(k, n)$ that W will be obtained k times is

$$P(k, n) = \binom{n}{k} p^k q^{n-k}. \quad (6.6)$$

This probability function is called the *Bernoulli* or, also, the *binomial* probability.

6-3 The random walk

There is another interesting problem in which the idea of probability is required. It is the problem of the “random walk.” In its simplest version, we imagine a “game” in which a “player” starts at the point $x = 0$ and at each “move” is required to take a step *either* forward (toward $+x$) *or* backward (toward $-x$). The choice is to be made *randomly*, determined, for example, by the toss of a coin. How shall we describe the resulting motion? In its general form the problem is related to the motion of atoms (or other particles) in a gas—called Brownian motion—and also to the combination of errors in measurements. You will see that the random-walk problem is closely related to the coin-tossing problem we have already discussed.

First, let us look at a few examples of a random walk. We may characterize the walker’s progress by the net distance D_N traveled in N steps. We show in the graph of Fig. 6-5 three examples of the path of a random walker. (We have used for the random sequence of choices the results of the coin tosses shown in Fig. 6-1.)

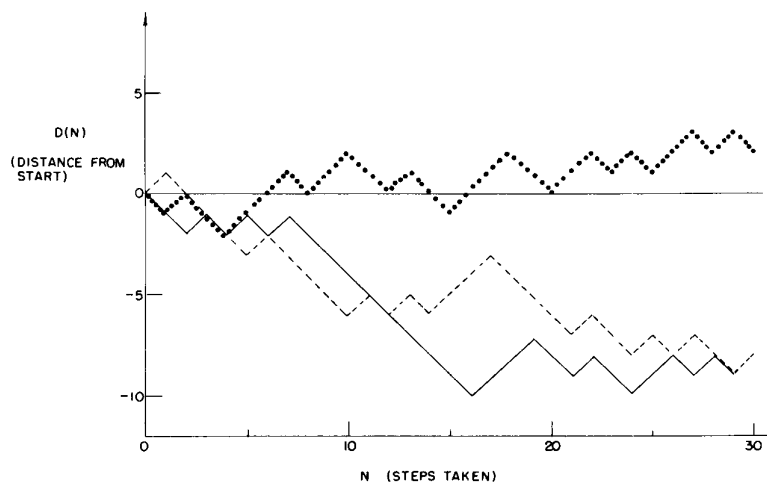


Fig. 6-5. The progress made in a random walk. The horizontal coordinate N is the total number of steps taken; the vertical coordinate $D(N)$ is the net distance moved from the starting position.

What can we say about such a motion? We might first ask: “How far does he get on the average?” We must *expect* that his average progress will be zero, since he is equally likely to go either forward or backward. But we have the feeling that as N increases, he is more likely to have strayed farther from the starting point. We might, therefore, ask what is his average distance travelled in *absolute value*, that is, what is the average of $|D|$. It is, however, more convenient to deal with another measure of “progress,” the square of the distance: D^2 is positive for either positive or negative motion, and is therefore a reasonable *measure* of such random wandering.

We can show that the expected value of D_N^2 is just N , the number of steps taken. By “expected value” we mean the probable value (our best guess), which we can think of as the *expected* average behavior in *many repeated* sequences. We represent such an expected value by $\langle D_N^2 \rangle$, and may refer to it also as the “mean square distance.” After one step, D^2 is always $+1$, so we have certainly $\langle D_1^2 \rangle = 1$. (All distances will be measured in terms of a unit of one step. We shall not continue to write the units of distance.)

The expected value of D_N^2 for $N > 1$ can be obtained from D_{N-1} . If, after $(N - 1)$ steps, we have D_{N-1} , then after N steps we have $D_N = D_{N-1} + 1$ or $D_N = D_{N-1} - 1$. For the squares,

$$D_N^2 = \begin{cases} D_{N-1}^2 + 2D_{N-1} + 1, \\ \text{or} \\ D_{N-1}^2 - 2D_{N-1} + 1. \end{cases} \quad (6.7)$$

In a number of independent sequences, we expect to obtain each value one-half of the time, so our average expectation is just the average of the two possible values. The expected value of D_N^2 is then $D_{N-1}^2 + 1$. In general, we should expect for D_{N-1}^2 its “expected value” $\langle D_{N-1}^2 \rangle$ (by definition!). So

$$\langle D_N^2 \rangle = \langle D_{N-1}^2 \rangle + 1. \quad (6.8)$$

We have already shown that $\langle D_1^2 \rangle = 1$; it follows then that

$$D_N^2 = N, \quad (6.9)$$

a particularly simple result!

If we wish a number like a distance, rather than a distance squared, to represent the “progress made away from the origin” in a random walk, we can use the “root-mean-square distance” D_{rms} :

$$D_{\text{rms}} = \sqrt{\langle D^2 \rangle} = \sqrt{N}. \quad (6.10)$$

We have pointed out that the random walk is closely similar in its mathematics to the coin-tossing game we considered at the beginning of the chapter. If we imagine the direction of each step to be in correspondence with the appearance of heads or tails in a coin toss, then D is just $N_H - N_T$, the difference in the number of heads and tails. Since $N_H + N_T = N$, the total number of steps (and tosses), we have $D = 2N_H - N$. We have derived earlier an expression for the expected distribution of N_H (also called k) and obtained the result of Eq. (6.5). Since N is just a constant, we have the corresponding distribution for D . (Since for every head more than $N/2$ there is a tail “missing,” we have the factor of 2 between N_H and D .) The graphs of Fig. 6-2 represent the distribution of distances we might get in 30 random steps (where $k = 15$ is to be read $D = 0$; $k = 16$, $D = 2$; etc.).

The variation of N_H from its expected value $N/2$ is

$$N_H - \frac{N}{2} = \frac{D}{2}. \quad (6.11)$$

The rms deviation is

$$\left(N_H - \frac{N}{2} \right)_{\text{rms}} = \frac{1}{2} \sqrt{N}. \quad (6.12)$$

According to our result for D_{rms} , we expect that the “typical” distance in 30 steps ought to be $\sqrt{30} = 5.5$, or a typical k should be about $5.5/2 = 2.8$ units from 15. We see that the “width” of the curve in Fig. 6-2, measured from the center, is just about 3 units, in agreement with this result.

We are now in a position to consider a question we have avoided until now. How shall we tell whether a coin is “honest” or “loaded”? We can give now at least a partial answer. For an honest coin, we expect the fraction of the times heads appears to be 0.5, that is,

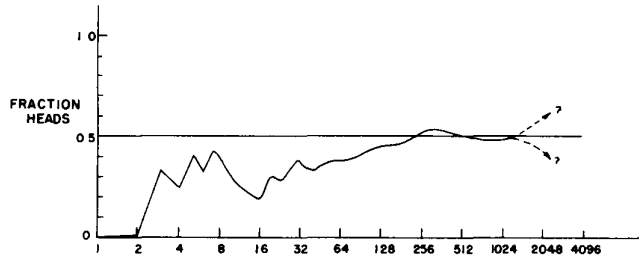
$$\frac{\langle N_H \rangle}{N} = 0.5. \quad (6.13)$$

We also expect an actual N_H to deviate from $N/2$ by about $\sqrt{N}/2$, or the fraction to deviate by

$$\frac{1}{N} \frac{\sqrt{N}}{2} = \frac{1}{2\sqrt{N}}.$$

The larger N is, the closer we expect the fraction N_H/N to be to one-half.

Fig. 6-6. The fraction of the tosses that gave heads in a particular sequence of N tosses of a penny.



In Fig. 6-6 we have plotted the fraction N_H/N for the coin tosses reported earlier in this chapter. We see the tendency for the fraction of heads to approach 0.5 for large N . Unfortunately, for any given run or combination of runs there is no *guarantee* that the observed deviation will be even *near* the *expected* deviation. There is always the finite chance that a large fluctuation—a long string of heads or tails—will give an arbitrarily large deviation. All we can say is that *if* the deviation is near the expected $1/2\sqrt{N}$ (say within a factor of 2 or 3), we have no reason to suspect the honesty of the coin. If it is much larger, we may be suspicious, but cannot prove, that the coin is loaded (or that the tosser is clever!).

We have also not considered how we should treat the case of a “coin” or some similar “chancy” object (say a stone that always lands in either of two positions) that we have good reason to believe *should* have a different probability for heads and tails. We have defined $P(H) = \langle N_H \rangle / N$. How shall we know what to *expect* for N_H ? In some cases, the best we can do is to observe the number of heads obtained in large numbers of tosses. For want of anything better, we must set $\langle N_H \rangle = N_H(\text{observed})$. (How could we expect anything else?) We must understand, however, that in such a case a different experiment, or a different observer, might conclude that $P(H)$ was different. We would *expect*, however, that the various answers should agree within the deviation $1/2\sqrt{N}$ [if $P(H)$ is near one-half]. An experimental physicist usually says that an “experimentally determined” probability has an “error,” and writes

$$P(H) = \frac{N_H}{N} \pm \frac{1}{2\sqrt{N}}. \quad (6.14)$$

There is an implication in such an expression that there *is* a “true” or “correct” probability which *could* be computed if we knew enough, and that the observation may be in “error” due to a fluctuation. There is, however, no way to make such thinking logically consistent. It is probably better to realize that the probability concept is in a sense subjective, that it is always based on uncertain knowledge, and that its quantitative evaluation is subject to change as we obtain more information.

6-4 A probability distribution

Let us return now to the random walk and consider a modification of it. Suppose that in addition to a random choice of the *direction* (+ or −) of each step, the *length* of each step also varied in some unpredictable way, the only condition being that *on the average* the step length was one unit. This case is more representative of something like the thermal motion of a molecule in a gas. If we call the length of a step S , then S may have any value at all, but most often will be “near” 1. To be specific, we shall let $\langle S^2 \rangle = 1$ or, equivalently, $S_{\text{rms}} = 1$. Our derivation for $\langle D^2 \rangle$ would proceed as before except that Eq. (6.8) would be changed now to read

$$\langle D_N^2 \rangle = \langle D_{N-1}^2 \rangle + \langle S^2 \rangle = \langle D_{N-1}^2 \rangle + 1. \quad (6.15)$$

We have, as before, that

$$\langle D_N^2 \rangle = N. \quad (6.16)$$

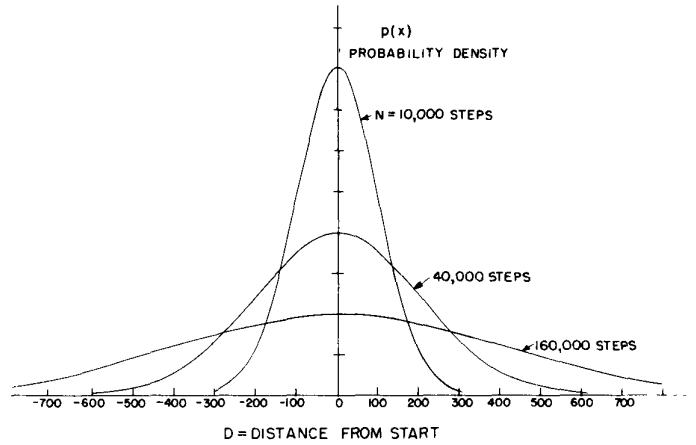


Fig. 6-7. The probability density for ending up at the distance D from the starting place in a random walk of N steps. (D is measured in units of the rms step length.)

What would we expect now for the distribution of distances D ? What is, for example, the probability that $D = 0$ after 30 steps? The answer is zero! The probability is zero that D will be *any particular* value, since there is no chance at all that the sum of the backward steps (of varying lengths) would exactly equal the sum of forward steps. We cannot plot a graph like that of Fig. 6-2.

We can, however, obtain a representation similar to that of Fig. 6-2, if we ask, not what is the probability of obtaining D exactly equal to 0, 1, or 2, but instead what is the probability of obtaining D *near* 0, 1, or 2. Let us define $P(x, \Delta x)$ as the probability that D will lie in the interval Δx located at x (say from x to $x + \Delta x$). We expect that for small Δx the chance of D landing in the interval is proportional to Δx , the width of the interval. So we can write

$$P(x, \Delta x) = p(x) \Delta x. \quad (6.17)$$

The function $p(x)$ is called the *probability density*.

The form of $p(x)$ will depend on N , the number of steps taken, and also on the distribution of individual step lengths. We cannot demonstrate the proofs here, but for large N , $p(x)$ is the *same* for all reasonable distributions in individual step lengths, and depends only on N . We plot $p(x)$ for three values of N in Fig. 6-7. You will notice that the “half-widths” (typical spread from $x = 0$) of these curves is \sqrt{N} , as we have shown it should be.

You may notice also that the value of $p(x)$ near zero is inversely proportional to \sqrt{N} . This comes about because the curves are all of a similar shape and their areas under the curves must all be equal. Since $p(x) \Delta x$ is the probability of finding D in Δx when Δx is small, we can determine the chance of finding D *somewhere* inside an arbitrary interval from x_1 to x_2 , by cutting the interval in a number of small increments Δx and evaluating the sum of the terms $p(x) \Delta x$ for each increment. The probability that D lands somewhere between x_1 and x_2 , which we may write $P(x_1 < D < x_2)$, is equal to the shaded area in Fig. 6-8. The smaller we take the increments Δx , the more correct is our result. We can write, therefore,

$$P(x_1 < D < x_2) = \sum p(x) \Delta x = \int_{x_1}^{x_2} p(x) dx. \quad (6.18)$$

The area under the whole curve is the probability that D lands somewhere (that is, has *some* value between $x = -\infty$ and $x = +\infty$). That probability is surely 1. We must have that

$$\int_{-\infty}^{+\infty} p(x) dx = 1. \quad (6.19)$$

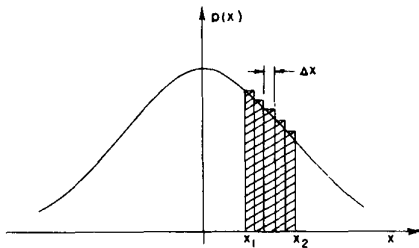


Fig. 6-8. The probability that the distance D traveled in a random walk is between x_1 and x_2 is the area under the curve of $p(x)$ from x_1 to x_2 .

Since the curves in Fig. 6-7 get wider in proportion to \sqrt{N} , their heights must be proportional to $1/\sqrt{N}$ to maintain the total area equal to 1.

The probability density function we have been describing is one that is encountered most commonly. It is known as the *normal* or *gaussian* probability density. It has the mathematical form

$$p(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-x^2/2\sigma^2}, \quad (6.20)$$

where σ is called the *standard deviation* and is given, in our case, by $\sigma = \sqrt{N}$ or, if the rms step size is different from 1, by $\sigma = \sqrt{N}S_{\text{rms}}$.

We remarked earlier that the motion of a molecule, or of any particle, in a gas is like a random walk. Suppose we open a bottle of an organic compound and let some of its vapor escape into the air. If there are air currents, so that the air is circulating, the currents will also carry the vapor with them. But even in *perfectly still air*, the vapor will gradually spread out—will diffuse—until it has penetrated throughout the room. We might detect it by its color or odor. The individual molecules of the organic vapor spread out in still air because of the molecular motions caused by collisions with other molecules. If we know the average “step” size, and the number of steps taken per second, we can find the probability that one, or several, molecules will be found at some distance from their starting point after any particular passage of time. As time passes, more steps are taken and the gas spreads out as in the successive curves of Fig. 6-7. In a later chapter, we shall find out how the step sizes and step frequencies are related to the temperature and pressure of a gas.

Earlier, we said that the pressure of a gas is due to the molecules bouncing against the walls of the container. When we come later to make a more quantitative description, we will wish to know how fast the molecules are going when they bounce, since the impact they make will depend on that speed. We cannot, however, speak of *the* speed of the molecules. It is necessary to use a probability description. A molecule may have any speed, but some speeds are more likely than others. We describe what is going on by saying that the probability that any particular molecule will have a speed between v and $v + \Delta v$ is $p(v) \Delta v$, where $p(v)$, a probability density, is a given function of the speed v . We shall see later how Maxwell, using common sense and the ideas of probability, was able to find a mathematical expression for $p(v)$. The form* of the function $p(v)$ is shown in Fig. 6-9. Velocities may have any value, but are most likely to be near the most probable or expected value $\langle v \rangle$.

We often think of the curve of Fig. 6-9 in a somewhat different way. If we consider the molecules in a typical container (with a volume of, say, one liter), then there are a very large number N of molecules present ($N \approx 10^{22}$). Since $p(v) \Delta v$ is the probability that *one* molecule will have its velocity in Δv , by our definition of probability we mean that the *expected* number $\langle \Delta N \rangle$ to be found with a velocity in the interval Δv is given by

$$\langle \Delta N \rangle = Np(v) \Delta v. \quad (6.21)$$

We call $Np(v)$ the “distribution in velocity.” The area under the curve between two velocities v_1 and v_2 , for example the shaded area in Fig. 6-9, represents [for the curve $Np(v)$] the expected number of molecules with velocities between v_1 and v_2 . Since with a gas we are usually dealing with large numbers of molecules, we expect the deviations from the expected numbers to be small (like $1/\sqrt{N}$), so we often neglect to say the “expected” number, and say instead: “The number of molecules with velocities between v_1 and v_2 is the area under the curve.” We should remember, however, that such statements are always about *probable* numbers.

* Maxwell's expression is $p(v) = Cv^2e^{-av^2}$, where a is a constant related to the temperature and C is chosen so that the total probability is one.

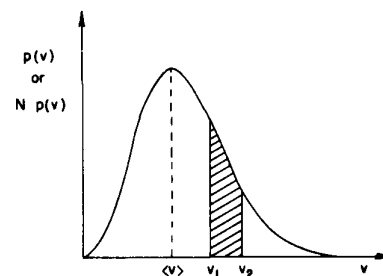


Fig. 6-9. The distribution of velocities of the molecules in a gas.

6-5 The uncertainty principle

The ideas of probability are certainly useful in describing the behavior of the 10^{22} or so molecules in a sample of a gas, for it is clearly impractical even to attempt to write down the position or velocity of each molecule. When probability was first applied to such problems, it was considered to be a *convenience*—a way of dealing with very complex situations. We now believe that the ideas of probability are *essential* to a description of atomic happenings. According to quantum mechanics, the mathematical theory of particles, there is always some uncertainty in the *specification* of positions and velocities. We can, at best, say that there is a certain probability that any particle will have a position near some coordinate x .

We can give a probability density $p_1(x)$, such that $p_1(x) \Delta x$ is the probability that the particle will be found between x and $x + \Delta x$. If the particle is reasonably well localized, say near x_1 , the function $p_1(x)$ might be given by the graph of Fig. 6-10(a). Similarly, we must specify the velocity of the particle by means of a probability density $p_2(v)$, with $p_2(v) \Delta v$ the probability that the velocity will be found between v and $v + \Delta v$.

It is one of the fundamental results of quantum mechanics that the two functions $p_1(x)$ and $p_2(v)$ cannot be chosen independently and, in particular, cannot both be made arbitrarily narrow. If we call the typical “width” of the $p_1(x)$ curve $[\Delta x]$, and that of the $p_2(v)$ curve $[\Delta v]$ (as shown in the figure), nature demands that the *product* of the two widths be at least as big as the number h/m , where m is the mass of the particle and h is a fundamental physical constant called *Planck’s constant*. We may write this basic relationship as

$$[\Delta x] \cdot [\Delta v] \geq h/m. \quad (6.22)$$

This equation is a statement of the *Heisenberg uncertainty principle* that we mentioned earlier.

Since the right-hand side of Eq. (6.22) is a constant, this equation says that if we try to “pin down” a particle by forcing it to be at a particular place, it ends up by having a high speed. Or if we try to force it to go very slowly, or at a precise velocity, it “spreads out” so that we do not know very well just where it is. Particles behave in a funny way!

The uncertainty principle describes an inherent fuzziness that must exist in any attempt to describe nature. Our most precise description of nature *must* be in terms of *probabilities*. There are some people who do not like this way of describing nature. They feel somehow that if they could only tell what is *really* going on with a particle, they could know its speed and position simultaneously. In the early days of the development of quantum mechanics, Einstein was quite worried about this problem. He used to shake his head and say, “But, surely God does not throw dice in determining how electrons should go!” He worried about that problem for a long time and he probably never really reconciled himself to the fact that this is the best description of nature that one can give. There are still one or two physicists who are working on the problem who have an intuitive conviction that it is possible somehow to describe the world in a different way and that all of this uncertainty about the way things are can be removed. No one has yet been successful.

The necessary uncertainty in our specification of the position of a particle becomes most important when we wish to describe the structure of atoms. In the hydrogen atom, which has a nucleus of one proton with one electron outside of the nucleus, the uncertainty in the position of the electron is as large as the atom itself! We cannot, therefore, properly speak of the electron moving in some “orbit” around the proton. The most we can say is that there is a certain *chance* $p(r) \Delta V$, of observing the electron in an element of volume ΔV at the distance r from the proton. The probability density $p(r)$ is given by quantum mechanics. For an undisturbed hydrogen atom $p(r) = Ae^{-r^2/a^2}$, which is a bell-shaped function like that in Fig. 6-8. The number a is the “typical” radius, where the function is decreasing rapidly. Since there is a small probability of finding the electron at distances

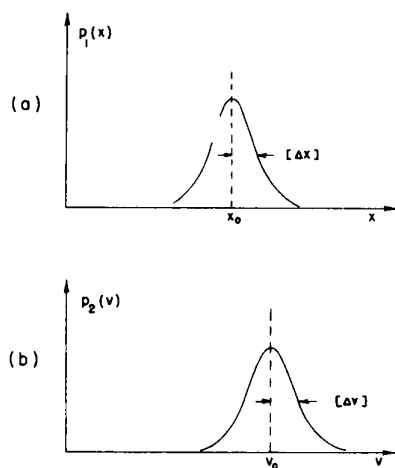
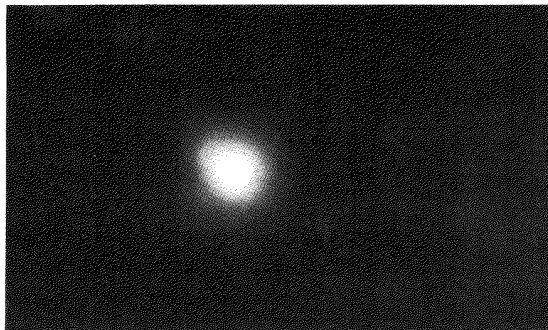


Fig. 6-10. Probability densities for observation of the position and velocity of a particle.

Fig. 6-11. A way of visualizing a hydrogen atom. The density (whiteness) of the cloud represents the probability density for observing the electron.



from the nucleus much greater than a , we may think of a as “the radius of the atom,” about 10^{-10} meter.

We can form an image of the hydrogen atom by imagining a “cloud” whose density is proportional to the probability density for observing the electron. A sample of such a cloud is shown in Fig. 6-11. Thus our best “picture” of a hydrogen atom is a nucleus surrounded by an “electron cloud” (although we *really* mean a “probability cloud”). The electron is there somewhere, but nature permits us to know only the *chance* of finding it at any particular place.

In its efforts to learn as much as possible about nature, modern physics has found that certain things can never be “known” with certainty. Much of our knowledge must always remain uncertain. The *most* we can know is in terms of probabilities.

The Theory of Gravitation

7-1 Planetary motions

In this chapter we shall discuss one of the most far-reaching generalizations of the human mind. While we are admiring the human mind, we should take some time off to stand in awe of a *nature* that could follow with such completeness and generality such an elegantly simple principle as the law of gravitation. What is this law of gravitation? It is that every object in the universe attracts every other object with a force which for any two bodies is proportional to the mass of each and varies inversely as the square of the distance between them. This statement can be expressed mathematically by the equation

$$F = G \frac{mm'}{r^2}.$$

If to this we add the fact that an object responds to a force by accelerating in the direction of the force by an amount that is inversely proportional to the mass of the object, we shall have said everything required, for a sufficiently talented mathematician could then deduce all the consequences of these two principles. However, since you are not assumed to be sufficiently talented yet, we shall discuss the consequences in more detail, and not just leave you with only these two bare principles. We shall briefly relate the story of the discovery of the law of gravitation and discuss some of its consequences, its effects on history, the mysteries that such a law entails, and some refinements of the law made by Einstein; we shall also discuss the relationships of the law to the other laws of physics. All this cannot be done in one chapter, but these subjects will be treated in due time in subsequent chapters.

The story begins with the ancients observing the motions of planets among the stars, and finally deducing that they went around the sun, a fact that was rediscovered later by Copernicus. Exactly *how* the planets went around the sun, with exactly *what motion*, took a little more work to discover. In the beginning of the fifteenth century there were great debates as to whether they really went around the sun or not. Tycho Brahe had an idea that was different from anything proposed by the ancients: his idea was that these debates about the nature of the motions of the planets would best be resolved if the actual positions of the planets in the sky were measured sufficiently accurately. If measurement showed exactly how the planets moved, then perhaps it would be possible to establish one or another viewpoint. This was a tremendous idea—that to find something out, it is better to perform some careful experiments than to carry on deep philosophical arguments. Pursuing this idea, Tycho Brahe studied the positions of the planets for many years in his observatory on the island of Hven, near Copenhagen. He made voluminous tables, which were then studied by the mathematician Kepler, after Tycho's death. Kepler discovered from the data some very beautiful and remarkable, but simple, laws regarding planetary motion.

7-2 Kepler's laws

First of all, Kepler found that each planet goes around the sun in a curve called an *ellipse*, with the sun at a focus of the ellipse. An ellipse is not just an oval, but is a very specific and precise curve that can be obtained by using two tacks, one at each focus, a loop of string, and a pencil; more mathematically, it

7-1 Planetary motions

7-2 Kepler's laws

7-3 Development of dynamics

7-4 Newton's law of gravitation

7-5 Universal gravitation

7-6 Cavendish's experiment

7-7 What is gravity?

7-8 Gravity and relativity

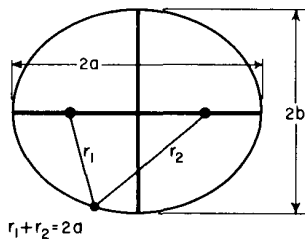


Fig. 7-1. An ellipse.

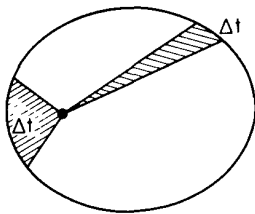


Fig. 7-2. Kepler's law of areas.

is the locus of all points the sum of whose distances from two fixed points (the foci) is a constant. Or, if you will, it is a foreshortened circle (Fig. 7-1).

Kepler's second observation was that the planets do not go around the sun at a uniform speed, but move faster when they are nearer the sun and more slowly when they are farther from the sun, in precisely this way: Suppose a planet is observed at any two successive times, let us say a week apart, and that the radius vector* is drawn to the planet for each observed position. The orbital arc traversed by the planet during the week, and the two radius vectors, bound a certain plane area, the shaded area shown in Fig. 7-2. If two similar observations are made a week apart, at a part of the orbit farther from the sun (where the planet moves more slowly), the similarly bounded area is exactly the same as in the first case. So, in accordance with the second law, the orbital speed of each planet is such that the radius "sweeps out" equal areas in equal times.

Finally, a third law was discovered by Kepler much later; this law is of a different category from the other two, because it deals not with only a single planet, but relates one planet to another. This law says that when the orbital period and orbit size of any two planets are compared, the periods are proportional to the $3/2$ power of the orbit size. In this statement the period is the time interval it takes a planet to go completely around its orbit, and the size is measured by the length of the greatest diameter of the elliptical orbit, technically known as the major axis. More simply, if the planets went in circles, as they nearly do, the time required to go around the circle would be proportional to the $3/2$ power of the diameter (or radius). Thus Kepler's three laws are:

- I. Each planet moves around the sun in an ellipse, with the sun at one focus.
- II. The radius vector from the sun to the planet sweeps out equal areas in equal intervals of time.
- III. The squares of the periods of any two planets are proportional to the cubes of the semimajor axes of their respective orbits: $T \sim a^{3/2}$.

7-3 Development of dynamics

While Kepler was discovering these laws, Galileo was studying the laws of motion. The problem was, what makes the planets go around? (In those days, one of the theories proposed was that the planets went around because behind them were invisible angels, beating their wings and driving the planets forward. You will see that this theory is now modified! It turns out that in order to keep the planets going around, the invisible angels must fly in a different direction and they have no wings. Otherwise, it is a somewhat similar theory!) Galileo discovered a very remarkable fact about motion, which was essential for understanding these laws. That is the principle of *inertia*—if something is moving, with nothing touching it and completely undisturbed, it will go on forever, coasting at a uniform speed in a straight line. (*Why* does it keep on coasting? We do not know, but that is the way it is.)

Newton modified this idea, saying that the only way to change the motion of a body is to use *force*. If the body speeds up, a force has been applied *in the direction of motion*. On the other hand, if its motion is changed to a new *direction*, a force has been applied *sideways*. Newton thus added the idea that a force is needed to change the speed *or the direction* of motion of a body. For example, if a stone is attached to a string and is whirling around in a circle, it takes a force to keep it in the circle. We have to *pull* on the string. In fact, the law is that the acceleration produced by the force is inversely proportional to the mass, or the force is proportional to the mass times the acceleration. The more massive a thing is, the stronger the force required to produce a given acceleration. (The mass can be measured by putting other stones on the end of the same string and making them go around the same circle at the same speed. In this way it is found that more or less force is required, the more massive object requiring more force.)

* A radius vector is a line drawn from the sun to any point in a planet's orbit.

The brilliant idea resulting from these considerations is that no *tangential* force is needed to keep a planet in its orbit (the angels do not have to fly tangentially) because the planet would coast in that direction anyway. If there were nothing at all to disturb it, the planet would go off in a *straight line*. But the actual motion deviates from the line on which the body would have gone if there were no force, the deviation being essentially *at right angles* to the motion, not in the direction of the motion. In other words, because of the principle of inertia, the force needed to control the motion of a planet around the sun is not a force *around* the sun but *toward* the sun. (If there is a force toward the sun, the sun might be the angel, of course!)

7-4 Newton's law of gravitation

From his better understanding of the theory of motion, Newton appreciated that *the sun* could be the seat or organization of forces that govern the motion of the planets. Newton proved to himself (and perhaps we shall be able to prove it soon) that the very fact that equal areas are swept out in equal times is a precise sign post of the proposition that all deviations are precisely *radial*—that the law of areas is a direct consequence of the idea that all of the forces are directed exactly *toward the sun*.

Next, by analyzing Kepler's third law it is possible to show that the farther away the planet, the weaker the forces. If two planets at different distances from the sun are compared, the analysis shows that the forces are inversely proportional to the squares of the respective distances. With the combination of the two laws, Newton concluded that there must be a force, inversely as the square of the distance, directed in a line between the two objects.

Being a man of considerable feeling for generalities, Newton supposed, of course, that this relationship applied more generally than just to the sun holding the planets. It was already known, for example, that the planet Jupiter had moons going around it as the moon of the earth goes around the earth, and Newton felt certain that each planet held its moons with a force. He already knew of the force holding *us* on the earth, so he proposed that this was a *universal* force—that *everything pulls everything else*.

The next problem was whether the pull of the earth on its people was the "same" as its pull on the moon, i.e., inversely as the square of the distance. If an object on the surface of the earth falls 16 feet in the first second after it is released from rest, how far does the moon fall in the same time? We might say that the moon does not fall at all. But if there were no force on the moon, it would go off in a straight line, whereas it goes in a circle instead, so it really *falls in* from where it would have been if there were no force at all. We can calculate from the radius of the moon's orbit (which is about 240,000 miles) and how long it takes to go around the earth (approximately 29 days), how far the moon moves in its orbit in 1 second, and can then calculate how far it falls in one second.* This distance turns out to be roughly 1/20 of an inch in a second. That fits very well with the inverse square law, because the earth's radius is 4000 miles, and if something which is 4000 miles from the center of the earth falls 16 feet in a second, something 240,000 miles, or 60 times as far away, should fall only 1/3600 of 16 feet, which also is roughly 1/20 of an inch. Wishing to put this theory of gravitation to a test by similar calculations, Newton made his calculations very carefully and found a discrepancy so large that he regarded the theory as contradicted by facts, and did not publish his results. Six years later a new measurement of the size of the earth showed that the astronomers had been using an incorrect distance to the moon. When Newton heard of this, he made the calculation again, with the corrected figures, and obtained beautiful agreement.

This idea that the moon "falls" is somewhat confusing, because, as you see, it does not come any *closer*. The idea is sufficiently interesting to merit further

* That is, how far the circle of the moon's orbit falls below the straight line tangent to it at the point where the moon was one second before.

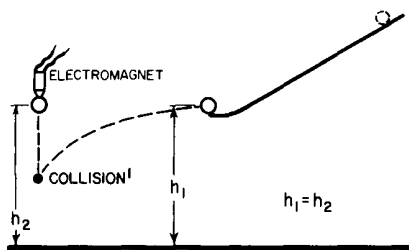


Fig. 7-3. Apparatus for showing the independence of vertical and horizontal motions.

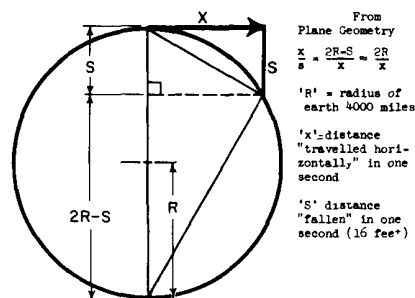


Fig. 7-4. Acceleration toward the center of a circular path. From plane geometry, $x/s = (2R - S)/x \approx 2R/x$, where R is the radius of the earth, 4000 miles; x is the distance "travelled horizontally" in one second; and S is the distance "fallen" in one second (16 feet).

explanation: the moon falls in the sense that *it falls away from the straight line that it would pursue if there were no forces*. Let us take an example on the surface of the earth. An object released near the earth's surface will fall 16 feet in the first second. An object shot out *horizontally* will also fall 16 feet; even though it is moving horizontally, it still falls the same 16 feet in the same time. Figure 7-3 shows an apparatus which demonstrates this. On the horizontal track is a ball which is going to be driven forward a little distance away. At the same height is a ball which is going to fall vertically, and there is an electrical switch arranged so that at the moment the first ball leaves the track, the second ball is released. That they come to the same depth at the same time is witnessed by the fact that they collide in midair. An object like a bullet, shot horizontally, might go a long way in one second—perhaps 2000 feet—but it will still fall 16 feet if it is aimed horizontally. What happens if we shoot a bullet faster and faster? Do not forget that the earth's surface is curved. If we shoot it fast enough, then when it falls 16 feet it may be at just the same height above the ground as it was before. How can that be? It still falls, but the earth curves away, so it falls "around" the earth. The question is, how far does it have to go in one second so that the earth is 16 feet below the horizon? In Fig. 7-4 we see the earth with its 4000-mile radius, and the tangential, straightline path that the bullet would take if there were no force. Now, if we use one of those wonderful theorems in geometry, which says that our tangent is the mean proportional between the two parts of the diameter cut by an equal chord, we see that the horizontal distance travelled is the mean proportional between the 16 feet fallen and the 8000-mile diameter of the earth. The square root of $(16/5280) \times 8000$ comes out very close to 5 miles. Thus we see that if the bullet moves at 5 miles a second, it then will continue to fall toward the earth at the same rate of 16 feet each second, but will never get any closer because the earth keeps curving away from it. Thus it was that Mr. Gagarin maintained himself in space while going 25,000 miles around the earth at approximately 5 miles per second. (He took a little longer because he was a little higher.)

Any great discovery of a new law is useful only if we can take more out than we put in. Now, Newton *used* the second and third of Kepler's laws to deduce his law of gravitation. What did he *predict*? First, his analysis of the moon's motion was a prediction because it connected the falling of objects on the earth's surface with that of the moon. Second, the question is, *is the orbit an ellipse*? We shall see in a later chapter how it is possible to calculate the motion exactly, and indeed one can prove that it should be an ellipse,* so no extra fact is needed to explain Kepler's *first* law. Thus Newton made his first powerful prediction.

The law of gravitation explains many phenomena not previously understood. For example, the pull of the moon on the earth causes the tides, hitherto mysterious. The moon pulls the water up under it and makes the tides—people had thought of that before, but they were not as clever as Newton, and so they thought there ought to be only one tide during the day. The reasoning was that the moon pulls the water up under it, making a high tide and a low tide, and since the earth spins underneath, that makes the tide at one station go up and down every 24 hours. Actually the tide goes up and down in 12 hours. Another school of thought claimed that the high tide should be on the other side of the earth because, so they argued, the moon pulls the earth away from the water! Both of these theories are wrong. It actually works like this: the pull of the moon for the earth and for the water is "balanced" at the center. But the water which is closer to the moon is pulled *more* than the average and the water which is farther away from it is pulled *less* than the average. Furthermore, the water can flow while the more rigid earth cannot. The true picture is a combination of these two things.

What do we mean by "balanced"? What balances? If the moon pulls the whole earth toward it, why doesn't the earth fall right "up" to the moon? Because the earth does the same trick as the moon, it goes in a circle around a point which is inside the earth but not at its center. The moon does not just go around the

* The proof is not given in this course.

earth, the earth and the moon both go around a central position, each falling toward this common position, as shown in Fig. 7-5. This motion around the common center is what balances the fall of each. So the earth is not going in a straight line either; it travels in a circle. The water on the far side is "unbalanced" because the moon's attraction there is weaker than it is at the center of the earth, where it just balances the "centrifugal force." The result of this imbalance is that the water rises up, away from the center of the earth. On the near side, the attraction from the moon is stronger, and the imbalance is in the opposite direction in space, but again *away* from the center of the earth. The net result is that we get *two* tidal bulges.

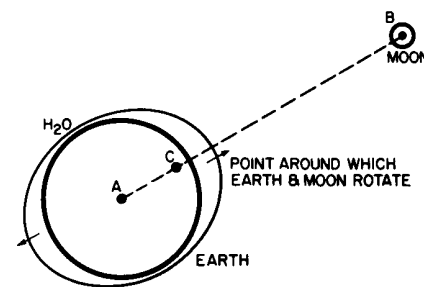


Fig. 7-5. The earth-moon system, with tides.

7-5 Universal gravitation

What else can we understand when we understand gravity? Everyone knows the earth is round. Why is the earth round? That is easy; it is due to gravitation. The earth can be understood to be round merely because everything attracts everything else and so it has attracted itself together as far as it can! If we go even further, the earth is not *exactly* a sphere because it is rotating, and this brings in centrifugal effects which tend to oppose gravity near the equator. It turns out that the earth should be elliptical, and we even get the right shape for the ellipse. We can thus deduce that the sun, the moon, and the earth should be (nearly) spheres, just from the law of gravitation.

What else can you do with the law of gravitation? If we look at the moons of Jupiter we can understand everything about the way they move around that planet. Incidentally, there was once a certain difficulty with the moons of Jupiter that is worth remarking on. These satellites were studied very carefully by Roemer, who noticed that the moons sometimes seemed to be ahead of schedule, and sometimes behind. (One can find their schedules by waiting a very long time and finding out how long it takes on the average for the moons to go around.) Now they were *ahead* when Jupiter was particularly *close* to the earth and they were *behind* when Jupiter was *farther* from the earth. This would have been a very difficult thing to explain according to the law of gravitation—it would have been, in fact, the death of this wonderful theory if there were no other explanation. If a law does not work even in *one place* where it ought to, it is just wrong. But the reason for this discrepancy was very simple and beautiful: it takes a little while to *see* the moons of Jupiter because of the time it takes light to travel from Jupiter to the earth. When Jupiter is closer to the earth the time is a little less, and when it is farther from the earth, the time is more. This is why moons appear to be, on the average, a little ahead or a little behind, depending on whether they are closer to or farther from the earth. This phenomenon showed that light does not travel instantaneously, and furnished the first estimate of the speed of light. This was done in 1656.

If all of the planets push and pull on each other, the force which controls, let us say, Jupiter in going around the sun is not just the force from the sun; there is also a pull from, say, Saturn. This force is not really strong, since the sun is much more massive than Saturn, but there is *some* pull, so the orbit of Jupiter should not be a perfect ellipse, and it is not; it is slightly off, and "wobbles" around the correct elliptical orbit. Such a motion is a little more complicated. Attempts were made to analyze the motions of Jupiter, Saturn, and Uranus on the basis of the law of gravitation. The effects of each of these planets on each other were calculated to see whether or not the tiny deviations and irregularities in these motions could be completely understood from this one law. Lo and behold, for Jupiter and Saturn, all was well, but Uranus was "weird." It behaved in a very peculiar manner. It was not travelling in an exact ellipse, but that was understandable, because of the attractions of Jupiter and Saturn. But even if allowance were made for these attractions, Uranus *still* was not going right, so the laws of gravitation were in danger of being overturned, a possibility that could not be ruled out. Two men, Adams and Leverrier, in England and France, independently,

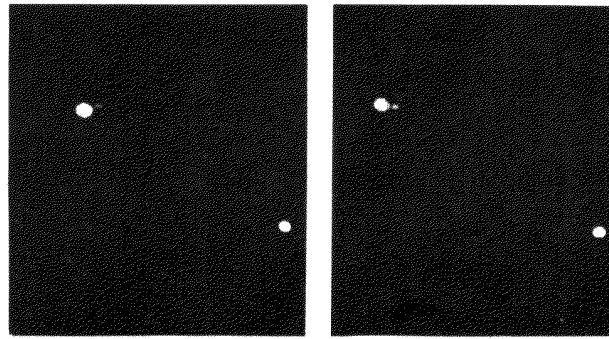


Fig. 7-6. A double-star system.

arrived at another possibility: perhaps there is *another* planet, dark and invisible, which men had not seen. This planet, *N*, could pull on Uranus. They calculated where such a planet would have to be in order to cause the observed perturbations. They sent messages to the respective observatories, saying, "Gentlemen, point your telescope to such and such a place, and you will see a new planet." It often depends on with whom you are working as to whether they pay any attention to you or not. They did pay attention to Leverrier; they looked, and there planet *N* was! The other observatory then also looked very quickly in the next few days and saw it too.

This discovery shows that Newton's laws are absolutely right in the solar system; but do they extend beyond the relatively small distances of the nearest planets? The first test lies in the question, do *stars* attract *each other* as well as planets? We have definite evidence that they do in the *double stars*. Figure 7-6 shows a double star—two stars very close together (there is also a third star in the picture so that we will know that the photograph was not turned). The stars are also shown as they appeared several years later. We see that, relative to the "fixed" star, the axis of the pair has rotated, i.e., the two stars are going around each other. Do they rotate according to Newton's laws? Careful measurements of the relative positions of one such double star system are shown in Fig. 7-7. There we see a beautiful ellipse, the measures starting in 1862 and going all the way around to 1904 (by now it must have gone around once more). Everything coincides with Newton's laws, except that the star Sirius A is *not at the focus*. Why should that be? Because the plane of the ellipse is not in the "plane of the sky." We are not looking at right angles to the orbit plane, and when an ellipse is viewed at a tilt, it remains an ellipse but the focus is no longer at the same place. Thus we can analyze double stars, moving about each other, according to the requirements of the gravitational law.

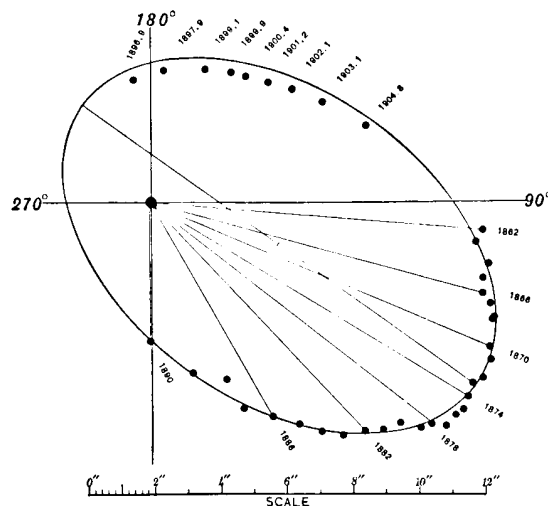


Fig. 7-7. Orbit of Sirius B with respect to Sirius A.

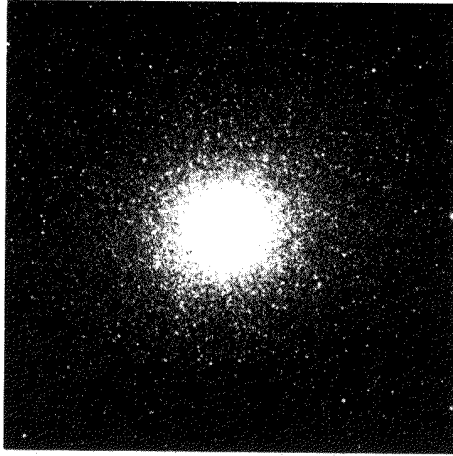


Fig. 7-8. A globular star cluster.

That the law of gravitation is true at even bigger distances is indicated in Fig. 7-8. If one cannot see gravitation acting here, he has no soul. This figure shows one of the most beautiful things in the sky—a globular star cluster. All of the dots are stars. Although they look as if they are packed solid toward the center, that is due to the fallibility of our instruments. Actually, the distances between even the centermost stars are very great and they very rarely collide. There are more stars in the interior than farther out, and as we move outward there are fewer and fewer. It is obvious that there is an attraction among these stars. It is clear that gravitation exists at these enormous dimensions, perhaps 100,000 times the size of the solar system. Let us now go further, and look at an *entire galaxy*, shown in Fig. 7-9. The shape of this galaxy indicates an obvious tendency for its matter to agglomerate. Of course we cannot prove that the law here is precisely inverse square, only that there is still an attraction, at this enormous dimension, that holds the whole thing together. One may say, “Well, that is all very clever but why is it not just a ball?” Because it is *spinning* and has *angular momentum* which it cannot give up as it contracts; it must contract mostly in a plane. (Incidentally, if you are looking for a good problem, the exact details of how the arms are formed and what determines the shapes of these galaxies has not been worked out.) It is, however, clear that the shape of the galaxy is due to gravitation even though the complexities of its structure have not yet allowed



Fig. 7-9. A galaxy.

us to analyze it completely. In a galaxy we have a scale of perhaps 50,000 to 100,000 light years. The earth's distance from the sun is $8\frac{1}{3}$ light *minutes*, so you can see how large these dimensions are.

Gravity appears to exist at even bigger dimensions, as indicated by Fig. 7-10, which shows many "little" things clustered together. This is a *cluster of galaxies*, just like a star cluster. Thus galaxies attract each other at such distances that they too are agglomerated into clusters. Perhaps gravitation exists even over distances of *tens of millions* of light years; so far as we now know, gravity seems to go out forever inversely as the square of the distance.

Not only can we understand the nebulae, but from the law of gravitation we can even get some ideas about the origin of the stars. If we have a big cloud of dust and gas, as indicated in Fig. 7-11, the gravitational attractions of the pieces of dust for one another might make them form little lumps. Barely visible in the figure are "little" black spots which may be the beginning of the accumulations of dust and gases which, due to their gravitation, begin to form stars. Whether we have ever seen a star form or not is still debatable. Figure 7-12 shows the one piece of evidence which suggests that we have. At the left is a picture of a region of gas with some stars in it taken in 1947, and at the right is another picture, taken only 7 years later, which shows two new bright spots. Has gas accumulated, has gravity acted hard enough and collected it into a ball big enough that the stellar nuclear reaction starts in the interior and turns it into a star? Perhaps, and perhaps not. It is unreasonable that in only seven years we should be so lucky as to see a star change itself into visible form; it is much less probable that we should see *two*!



Fig. 7-10. A cluster of galaxies.

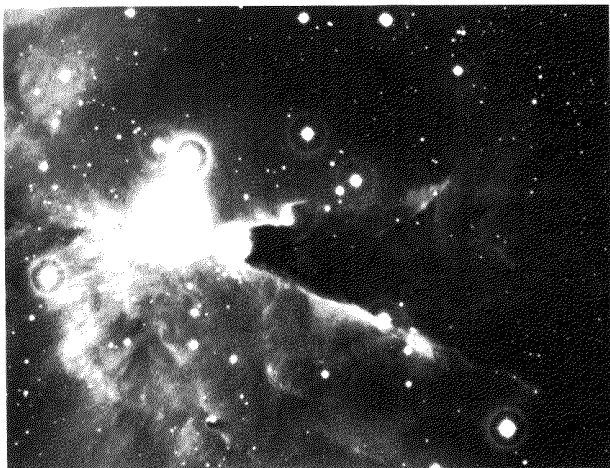


Fig. 7-11. An interstellar dust cloud.

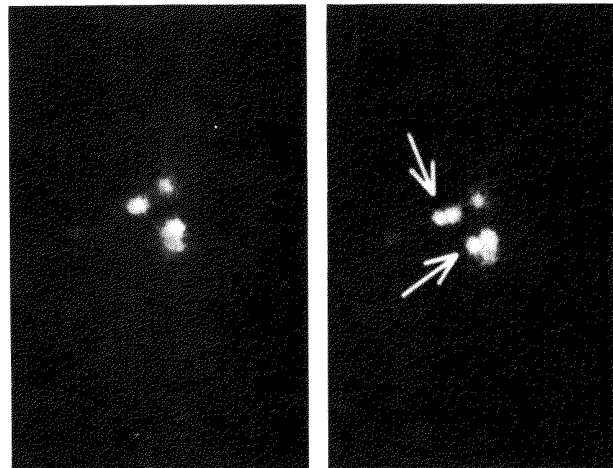


Fig. 7-12. The formation of new stars?

7-6 Cavendish's experiment

Gravitation, therefore, extends over enormous distances. But if there is a force between *any* pair of objects, we ought to be able to measure the force between our own objects. Instead of having to watch the stars go around each other, why can we not take a ball of lead and a marble and watch the marble go toward the ball of lead? The difficulty of this experiment when done in such a simple manner is the very weakness or delicacy of the force. It must be done with extreme care, which means covering the apparatus to keep the air out, making sure it is not electrically charged, and so on; then the force can be measured. It was first measured by Cavendish with an apparatus which is schematically indicated in Fig. 7-13. This first demonstrated the direct force between two large, fixed balls of lead and two smaller balls of lead on the ends of an arm supported by a very fine fiber, called a torsion fiber. By measuring how much the fiber gets twisted, one can measure the strength of the force, verify that it is inversely proportional to the square of the distance, and determine how strong it is. Thus, one may accurately determine the coefficient G in the formula

$$F = G \frac{mm'}{r^2}.$$

All the masses and distances are known. You say, "We knew it already for the earth." Yes, but we did not know the *mass* of the earth. By knowing G from this experiment and by knowing how strongly the earth attracts, we can indirectly learn how great is the mass of the earth! This experiment has been called "weighing the earth." Cavendish claimed he was weighing the earth, but what he was measuring was the coefficient G of the gravity law. This is the only way in which the mass of the earth can be determined. G turns out to be

$$6.670 \times 10^{-11} \text{ newton} \cdot \text{m}^2/\text{kg}^2.$$

It is hard to exaggerate the importance of the effect on the history of science produced by this great success of the theory of gravitation. Compare the confusion, the lack of confidence, the incomplete knowledge that prevailed in the earlier ages, when there were endless debates and paradoxes, with the clarity and simplicity of this law—this fact that all the moons and planets and stars have such a *simple rule* to govern them, and further that man could *understand* it and deduce how the planets should move! This is the reason for the success of the sciences in following years, for it gave hope that the other phenomena of the world might also have such beautifully simple laws.

7-7 What is gravity?

But is this such a simple law? What about the machinery of it? All we have done is to describe *how* the earth moves around the sun, but we have not said *what makes it go*. Newton made no hypotheses about this; he was satisfied to find *what* it did without getting into the machinery of it. *No one has since given any machinery*. It is characteristic of the physical laws that they have this abstract character. The law of conservation of energy is a theorem concerning quantities that have to be calculated and added together, with no mention of the machinery, and likewise the great laws of mechanics are quantitative mathematical laws for which no machinery is available. Why can we use mathematics to describe nature without a mechanism behind it? No one knows. We have to keep going because we find out more that way.

Many mechanisms for gravitation have been suggested. It is interesting to consider one of these, which many people have thought of from time to time. At first, one is quite excited and happy when he "discovers" it, but he soon finds that it is not correct. It was first discovered about 1750. Suppose there were many particles moving in space at a very high speed in all directions and being only slightly absorbed in going through matter. When they *are* absorbed, they give an impulse to the earth. However, since there are as many going one way as another, the

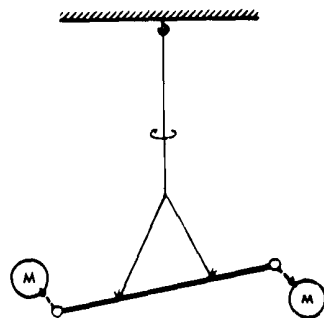


Fig. 7-13. A simplified diagram of the apparatus used by Cavendish to verify the law of universal gravitation for small objects and to measure the gravitational constant G .

turns out that if we consider the structure of the sun—the balance between the weight of its material and the rate at which radiant energy is generated inside it—we can deduce that if the gravity were 10 percent stronger, the sun would be much more than 10 percent brighter—by the *sixth power* of the gravity constant! If we calculate what happens to the orbit of the earth when the gravity is changing, we find that the earth was then *closer in*. Altogether, the earth would be about 100 degrees centigrade hotter, and all of the water would not have been in the sea, but vapor in the air, so life would not have started in the sea. So we do *not* now believe that the gravity constant is changing with the age of the universe. But such arguments as the one we have just given are not very convincing, and the subject is not completely closed.

It is a fact that the force of gravitation is proportional to the *mass*, the quantity which is fundamentally a measure of *inertia*—of how hard it is to hold something which is going around in a circle. Therefore two objects, one heavy and one light, going around a larger object in the same circle at the same speed because of gravity, will stay together because to go in a circle *requires* a force which is stronger for a bigger mass. That is, the gravity is stronger for a given mass in *just the right proportion* so that the two objects will go around together. If one object were inside the other it would *stay* inside; it is a perfect balance. Therefore, Gagarin or Titov would find things “weightless” inside a space ship; if they happened to let go of a piece of chalk, for example, it would go around the earth in exactly the same way as the whole space ship, and so it would appear to remain suspended before them in space. It is very interesting that this force is *exactly* proportional to the mass with great precision, because if it were not exactly proportional there would be some effect by which inertia and weight would differ. The absence of such an effect has been checked with great accuracy by an experiment done first by Eotvos in 1909 and more recently by Dicke. For all substances tried, the masses and weights are exactly proportional within 1 part in 1,000,000,000, or less. This is a remarkable experiment.

7-8 Gravity and relativity

Another topic deserving discussion is Einstein’s modification of Newton’s law of gravitation. In spite of all the excitement it created, Newton’s law of gravitation is not correct! It was modified by Einstein to take into account the theory of relativity. According to Newton, the gravitational effect is instantaneous, that is, if we were to move a mass, we would at once feel a new force because of the new position of that mass; by such means we could send signals at infinite speed. Einstein advanced arguments which suggest that we *cannot send signals faster than the speed of light*, so the law of gravitation must be wrong. By correcting it to take the delays into account, we have a new law, called Einstein’s law of gravitation. One feature of this new law which is quite easy to understand is this: In the Einstein relativity theory, anything which has *energy* has mass—mass in the sense that it is attracted gravitationally. Even light, which has an energy, has a “mass.” When a light beam, which has energy in it, comes past the sun there is an attraction on it by the sun. Thus the light does not go straight, but is deflected. During the eclipse of the sun, for example, the stars which are around the sun should appear displaced from where they would be if the sun were not there, and this has been observed.

Finally, let us compare gravitation with other theories. In recent years we have discovered that all mass is made of tiny particles and that there are several kinds of interactions, such as nuclear forces, etc. None of these nuclear or electrical forces has yet been found to explain gravitation. The quantum-mechanical aspects of nature have not yet been carried over to gravitation. When the scale is so small that we need the quantum effects, the gravitational effects are so weak that the need for a quantum theory of gravitation has not yet developed. On the other hand, for consistency in our physical theories it would be important to see whether Newton’s law modified to Einstein’s law can be further modified to be consistent with the uncertainty principle. This last modification has not yet been completed.

Motion

8-1 Description of motion

In order to find the laws governing the various changes that take place in bodies as time goes on, we must be able to *describe* the changes and have some way to record them. The simplest change to observe in a body is the apparent change in its position with time, which we call motion. Let us consider some solid object with a permanent mark, which we shall call a point, that we can observe. We shall discuss the motion of the little marker, which might be the radiator cap of an automobile or the center of a falling ball, and shall try to describe the fact that it moves and how it moves.

These examples may sound trivial, but many subtleties enter into the description of change. Some changes are more difficult to describe than the motion of a point on a solid object, for example the speed of drift of a cloud that is drifting very slowly, but rapidly forming or evaporating, or the change of a woman's mind. We do not know a simple way to analyze a change of mind, but since the cloud can be represented or described by many molecules, perhaps we can describe the motion of the cloud in principle by describing the motion of all its individual molecules. Likewise, perhaps even the changes in the mind may have a parallel in changes of the atoms inside the brain, but we have no such knowledge yet.

At any rate, that is why we begin with the motion of points; perhaps we should think of them as atoms, but it is probably better to be more rough in the beginning and simply to think of some kind of small objects—small, that is, compared with the distance moved. For instance, in describing the motion of a car that is going a hundred miles, we do not have to distinguish between the front and the back of the car. To be sure, there are slight differences, but for rough purposes we say “the car,” and likewise it does not matter that our points are not absolute points; for our present purposes it is not necessary to be extremely precise. Also, while we take a first look at this subject we are going to forget about the three dimensions of the world. We shall just concentrate on moving in one direction, as in a car on one road. We shall return to three dimensions after we see how to describe motion in one dimension. Now, you may say, “This is all some kind of trivia,” and indeed it is. How can we describe such a one-dimensional motion—let us say, of a car? Nothing could be simpler. Among many possible ways, one would be the following. To determine the position of the car at different times, we measure its distance from the starting point and record all the observations. In Table 8-1, s represents the distance of the car, in feet, from the starting point, and t represents the time in minutes. The first line in the table represents zero distance and zero time—the car has not started yet. After one minute it has started and has gone 1200 feet. Then in two minutes, it goes farther—notice that it picked up more distance in the second minute—it has accelerated; but something happened between 3 and 4 and even more so at 5—it stopped at a light perhaps? Then it speeds up again and goes 13,000 feet by the end of 6 minutes, 18,000 feet at the end of 7 minutes, and 23,500 feet in 8 minutes; at 9 minutes it has advanced to only 24,000 feet, because in the last minute it was stopped by a cop.

That is one way to describe the motion. Another way is by means of a graph. If we plot the time horizontally and the distance vertically, we obtain a curve something like that shown in Fig. 8-1. As the time increases, the distance increases, at first very slowly and then more rapidly, and very slowly again for a little while at 4 minutes; then it increases again for a few minutes and finally, at 9 minutes, appears to have stopped increasing. These observations can be made from the

8-1 Description of motion

8-2 Speed

8-3 Speed as a derivative

8-4 Distance as an integral

8-5 Acceleration

Table 8-1

t (min)	s (ft)
0	0
1	1200
2	4000
3	9000
4	9500
5	9600
6	13000
7	18000
8	23500
9	24000

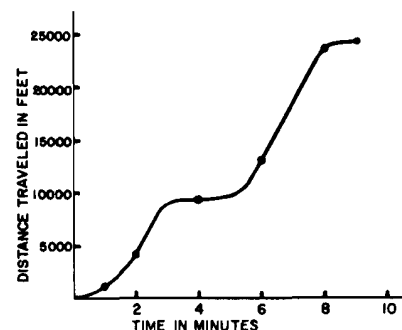


Fig. 8-1. Graph of distance versus time for the car.

Table 8-2

t (sec)	s (ft)
0	0
1	16
2	64
3	144
4	256
5	400
6	576

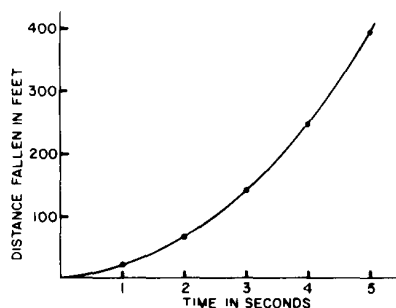


Fig. 8-2. Graph of distance versus time for a falling body.

graph, without a table. Obviously, for a complete description one would have to know where the car is at the half-minute marks, too, but we suppose that the graph means something, that the car has some position at all the intermediate times.

The motion of a car is complicated. For another example we take something that moves in a simpler manner, following more simple laws: a falling ball. Table 8-2 gives the time in seconds and the distance in feet for a falling body. At zero seconds the ball starts out at zero feet, and at the end of 1 second it has fallen 16 feet. At the end of 2 seconds, it has fallen 64 feet, at the end of 3 seconds, 144 feet, and so on; if the tabulated numbers are plotted, we get the nice parabolic curve shown in Fig. 8-2. The formula for this curve can be written as

$$s = 16t^2. \quad (8.1)$$

This formula enables us to calculate the distances at any time. You might say there ought to be a formula for the first graph too. Actually, one may write such a formula abstractly, as

$$s = f(t), \quad (8.2)$$

meaning that s is some quantity depending on t or, in mathematical phraseology, s is a function of t . Since we do not know what the function is, there is no way we can write it in definite algebraic form.

We have now seen two examples of motion, adequately described with very simple ideas, no subtleties. However, there *are* subtleties—several of them. In the first place, what do we mean by *time* and *space*? It turns out that these deep philosophical questions have to be analyzed very carefully in physics, and this is not so easy to do. The theory of relativity shows that our ideas of space and time are not as simple as one might think at first sight. However, for our present purposes, for the accuracy that we need at first, we need not be very careful about defining things precisely. Perhaps you say, “That’s a terrible thing—I learned that in science we have to define *everything* precisely.” We cannot define *anything* precisely! If we attempt to, we get into that paralysis of thought that comes to philosophers, who sit opposite each other, one saying to the other, “You don’t know what you are talking about!” The second one says, “What do you mean by *know*? What do you mean by *talking*? What do you mean by *you*?,” and so on. In order to be able to talk constructively, we just have to agree that we are talking about roughly the same thing. You know as much about time as we need for the present, but remember that there are some subtleties that have to be discussed; we shall discuss them later.

Another subtlety involved, and already mentioned, is that it should be possible to imagine that the moving point we are observing is always located somewhere. (Of course when we are looking at it, there it is, but maybe when we look away it isn’t there.) It turns out that in the motion of atoms, that idea also is false—we cannot find a marker on an atom and watch it move. That subtlety we shall have to get around in quantum mechanics. But we are first going to learn what the problems are before introducing the complications, and *then* we shall be in a better position to make corrections, in the light of the more recent knowledge of the subject. We shall, therefore, take a simple point of view about time and space. We know what these concepts are in a rough way, and those who have driven a car know what speed means.

8-2 Speed

Even though we know roughly what “speed” means, there are still some rather deep subtleties; consider that the learned Greeks were never able to adequately describe problems involving velocity. The subtlety comes when we try to comprehend exactly what is meant by “speed.” The Greeks got very confused about this, and a new branch of mathematics had to be discovered beyond the geometry and algebra of the Greeks, Arabs, and Babylonians. As an illustration of the difficulty, try to solve this problem by sheer algebra: A balloon is being inflated so

that the volume of the balloon is increasing at the rate of 100 cm^3 per second; at what speed is the radius increasing when the volume is 1000 cm^3 ? The Greeks were somewhat confused by such problems, being helped, of course, by some very confusing Greeks. To show that there were difficulties in reasoning about speed at the time, Zeno produced a large number of paradoxes, of which we shall mention one to illustrate his point that there are obvious difficulties in thinking about motion. "Listen," he says, "to the following argument: Achilles runs 10 times as fast as a tortoise, nevertheless he can never catch the tortoise. For, suppose that they start in a race where the tortoise is 100 meters ahead of Achilles; then when Achilles has run the 100 meters to the place where the tortoise was, the tortoise has proceeded 10 meters, having run one-tenth as fast. Now, Achilles has to run another 10 meters to catch up with the tortoise, but on arriving at the end of that run, he finds that the tortoise is still 1 meter ahead of him; running another meter, he finds the tortoise 10 centimeters ahead, and so on, *ad infinitum*. Therefore, at any moment the tortoise is always ahead of Achilles and Achilles can never catch up with the tortoise." What is wrong with that? It is that a finite amount of time can be divided into an infinite number of pieces, just as a length of line can be divided into an infinite number of pieces by dividing repeatedly by two. And so, although there are an infinite number of steps (in the argument) to the point at which Achilles reaches the tortoise, it doesn't mean that there is an infinite amount of *time*. We can see from this example that there are indeed some subtleties in reasoning about speed.

In order to get to the subtleties in a clearer fashion, we remind you of a joke which you surely must have heard. At the point where the lady in the car is caught by a cop, the cop comes up to her and says, "Lady, you were going 60 miles an hour!" She says, "That's impossible, sir, I was travelling for only seven minutes. It is ridiculous—how can I go 60 miles an hour when I wasn't going an hour?" How would you answer her if you were the cop? Of course, if you were really the cop, then no subtleties are involved; it is very simple: you say, "Tell that to the judge!" But let us suppose that we do not have that escape and we make a more honest, intellectual attack on the problem, and try to explain to this lady what we mean by the idea that she was going 60 miles an hour. Just what *do* we mean? We say, "What we mean, lady, is this: if you kept on going the same way as you are going now, in the next hour you would go 60 miles." She could say, "Well, my foot was off the accelerator and the car was slowing down, so if I kept on going that way it would not go 60 miles." Or consider the falling ball and suppose we want to know its speed at the time three seconds if the ball kept on going the way it is going. What does that mean—kept on *accelerating*, going faster? No—kept on going with the same *velocity*. But that is what we are trying to define! For if the ball keeps on going the way it is going, it will just keep on going the way it is going. Thus we need to define the velocity better. What has to be kept the same? The lady can also argue this way: "If I kept on going the way I'm going for one more hour, I would run into that wall at the end of the street!" It is not so easy to say what we mean.

Many physicists think that measurement is the only definition of anything. Obviously, then, we should use the instrument that measures the speed—the speedometer—and say, "Look, lady, your speedometer reads 60." So she says, "My speedometer is broken and didn't read at all." Does that mean the car is standing still? We believe that there is something to measure before we build the speedometer. Only then can we say, for example, "The speedometer isn't working right," or "the speedometer is broken." That would be a meaningless sentence if the velocity had no meaning independent of the speedometer. So we have in our minds, obviously, an idea that is independent of the speedometer, and the speedometer is meant only to measure this idea. So let us see if we can get a better definition of the idea. We say, "Yes, of course, before you went an hour, you would hit that wall, but if you went one second, you would go 88 feet; lady, you were going 88 feet per second, and if you kept on going, the next second it would be 88 feet, and the wall down there is farther away than that." She says, "Yes, but there's no law against going 88 feet per second! There is only a law

against going 60 miles an hour.” “But,” we reply, “it’s the same thing.” If it is the same thing, it should not be necessary to go into this circumlocution about 88 feet per second. In fact, the falling ball could not keep going the same way even one second because it would be changing speed, and we shall have to define speed somehow.

Now we seem to be getting on the right track; it goes something like this: If the lady kept on going for another $1/1000$ of an hour, she would go $1/1000$ of 60 miles. In other words, she does not have to keep on going for the whole hour; the point is that *for a moment* she is going at that speed. Now what that means is that if she went just a little bit more in time, the extra distance she goes would be the same as that of a car that goes at a *steady* speed of 60 miles an hour. Perhaps the idea of the 88 feet per second is right; we see how far she went in the last second, divide by 88 feet, and if it comes out 1 the speed was 60 miles an hour. In other words, we can find the speed in this way: We ask, how far do we go in a very short time? We divide that distance by the time, and that gives the speed. But the time should be made as short as possible, the shorter the better, because some change could take place during that time. If we take the time of a falling body as an hour, the idea is ridiculous. If we take it as a second, the result is pretty good for a car, because there is not much change in speed, but not for a falling body; so in order to get the speed more and more accurately, we should take a smaller and smaller time interval. What we should do is take a millionth of a second, and divide that distance by a millionth of a second. The result gives the distance per second, which is what we mean by the velocity, so we can define it that way. That is a successful answer for the lady, or rather, that is the definition that we are going to use.

The foregoing definition involves a new idea, an idea that was not available to the Greeks in a general form. That idea was to take an *infinitesimal distance* and the corresponding *infinitesimal time*, form the ratio, and watch what happens to that ratio as the time that we use gets smaller and smaller and smaller. In other words, take a limit of the distance travelled divided by the time required, as the time taken gets smaller and smaller, *ad infinitum*. This idea was invented by Newton and by Leibnitz, independently, and is the beginning of a new branch of mathematics, called the *differential calculus*. Calculus was invented in order to describe motion, and its first application was to the problem of defining what is meant by going “60 miles an hour.”

Let us try to define velocity a little better. Suppose that in a short time, ϵ , the car or other body goes a short distance x ; then the velocity, v , is defined as

$$v = x/\epsilon,$$

an approximation that becomes better and better as the ϵ is taken smaller and smaller. If a mathematical expression is desired, we can say that the velocity equals the limit as the ϵ is made to go smaller and smaller in the expression x/ϵ , or

$$v = \lim_{\epsilon \rightarrow 0} \frac{x}{\epsilon}. \quad (8.3)$$

We cannot do the same thing with the lady in the car, because the table is incomplete. We know only where she was at intervals of one minute; we can get a rough idea that she was going 5000 ft/min during the 7th minute, but we do not know, at exactly the moment 7 minutes, whether she had been speeding up and the speed was 4900 ft/min at the beginning of the 6th minute, and is now 5100 ft/min, or something else, because we do not have the exact details in between. So only if the table were completed with an infinite number of entries could we really calculate the velocity from such a table. On the other hand, when we have a complete mathematical formula, as in the case of a falling body (Eq. 8.1), then it is possible to calculate the velocity, because we can calculate the position at any time whatsoever.

Let us take as an example the problem of determining the velocity of the falling ball at the particular time 5 seconds. One way to do this is to see from

Table 8-2 what it did in the 5th second; it went $400 - 256 = 144$ ft, so it is going 144 ft/sec; however, that is wrong, because the speed is changing; *on the average* it is 144 ft/sec during this interval, but the ball is speeding up and is really going faster than 144 ft/sec. We want to find out *exactly how fast*. The technique involved in this process is the following: We know where the ball was at 5 sec. At 5.1 sec, the distance that it has gone all together is $16(5.1)^2 = 416.16$ ft (see Eq. 8.1). At 5 sec it had already fallen 400 ft; in the last tenth of a second it fell $416.16 - 400 = 16.16$ ft. Since 16.16 ft in 0.1 sec is the same as 161.6 ft/sec, that is the speed more or less, but it is not exactly correct. Is that the speed at 5, or at 5.1, or halfway between at 5.05 sec, or when *is* that the speed? Never mind—the problem was to find the speed *at 5 seconds*, and we do not have exactly that; we have to do a better job. So, we take one-thousandth of a second more than 5 sec, or 5.001 sec, and calculate the total fall as

$$s = 16(5.001)^2 = 16(25.010001) = 400.160016 \text{ ft.}$$

In the last 0.001 sec the ball fell 0.160016 ft, and if we divide this number by 0.001 sec we obtain the speed as 160.016 ft/sec. That is closer, very close, but it is *still not exact*. It should now be evident what we must do to find the speed exactly. To perform the mathematics we state the problem a little more abstractly: to find the velocity at a special time, t_0 , which in the original problem was 5 sec. Now the distance at t_0 , which we call s_0 , is $16t_0^2$, or 400 ft in this case. In order to find the velocity, we ask, “At the time $t_0 +$ (a little bit), or $t_0 + \epsilon$, where is the body?” The new position is $16(t_0 + \epsilon)^2 = 16t_0^2 + 32t_0\epsilon + 16\epsilon^2$. So it is farther along than it was before, because before it was only $16t_0^2$. This distance we shall call $s_0 +$ (a little bit more), or $s_0 + x$ (if x is the extra bit). Now if we subtract the distance at t_0 from the distance at $t_0 + \epsilon$, we get x , the extra distance gone, as $x = 32t_0 \cdot \epsilon + 16\epsilon^2$. Our first approximation to the velocity is

$$v = \frac{x}{\epsilon} = 32t_0 + 16\epsilon. \quad (8.4)$$

The true velocity is the value of this ratio, x/ϵ , when ϵ becomes vanishingly small. In other words, after forming the ratio, we take the limit as ϵ gets smaller and smaller, that is, approaches 0. The equation reduces to,

$$v \text{ (at time } t_0) = 32t_0.$$

In our problem, $t_0 = 5$ sec, so the solution is $v = 32 \times 5 = 160$ ft/sec. A few lines above, where we took ϵ as 0.1 and 0.01 sec successively, the value we got for v was a little more than this, but now we see that the actual velocity is precisely 160 ft/sec.

8-3 Speed as a derivative

The procedure we have just carried out is performed so often in mathematics that for convenience special notations have been assigned to our quantities ϵ and x . In this notation, the ϵ used above becomes Δt and x becomes Δs . This Δt means “an extra bit of t ,” and carries an implication that it can be made smaller. The prefix Δ is not a multiplier, any more than $\sin \theta$ means $s \cdot i \cdot n \cdot \theta$ —it simply defines a time increment, and reminds us of its special character. Δs has an analogous meaning for the distance s . Since Δ is not a factor, it cannot be cancelled in the ratio $\Delta s/\Delta t$ to give s/t , any more than the ratio $\sin \theta/\sin 2\theta$ can be reduced to $1/2$ by cancellation. In this notation, velocity is equal to the limit of $\Delta s/\Delta t$ when Δt gets smaller, or

$$v = \lim_{\Delta t \rightarrow 0} \frac{\Delta s}{\Delta t}. \quad (8.5)$$

This is really the same as our previous expression (8.3) with ϵ and x , but it has the advantage of showing that something is changing, and it keeps track of what is changing.

Incidentally, to a good approximation we have another law, which says that the change in distance of a moving point is the velocity times the time interval, or $\Delta s = v \Delta t$. This statement is true only if the velocity is not changing during that time interval, and this condition is true only in the limit as Δt goes to 0. Physicists like to write it $ds = v dt$, because by dt they mean Δt in circumstances in which it is very small; with this understanding, the expression is valid to a close approximation. If Δt is too long, the velocity might change during the interval, and the approximation would become less accurate. For a time dt , approaching zero, $ds = v dt$ precisely. In this notation we can write (8.5) as

$$v = \lim_{\Delta t \rightarrow 0} \frac{\Delta s}{\Delta t} = \frac{ds}{dt}.$$

The quantity ds/dt which we found above is called the “derivative of s with respect to t ” (this language helps to keep track of what was changed), and the complicated process of finding it is called finding a derivative, or differentiating. The ds ’s and dt ’s which appear separately are called *differentials*. To familiarize you with the words, we say we found the derivative of the function $16t^2$, or the derivative (with respect to t) of $16t^2$ is $32t$. When we get used to the words, the ideas are more easily understood. For practice, let us find the derivative of a more complicated function. We shall consider the formula $s = At^3 + Bt + C$, which might describe the motion of a point. The letters A , B , and C represent constant numbers, as in the familiar general form of a quadratic equation. Starting from the formula for the motion, we wish to find the velocity at any time. To find the velocity in the more elegant manner, we change t to $t + \Delta t$ and note that s is then changed to $s + \text{some } \Delta s$; then we find the Δs in terms of Δt . That is to say,

$$\begin{aligned} s + \Delta s &= A(t + \Delta t)^3 + B(t + \Delta t) + C \\ &= At^3 + Bt + C + 3At^2\Delta t + B\Delta t + 3At(\Delta t)^2 + A(\Delta t)^3, \end{aligned}$$

but since

$$s = At^3 + Bt + C,$$

we find that

$$\Delta s = 3At^2\Delta t + B\Delta t + 3At(\Delta t)^2 + A(\Delta t)^3.$$

But we do not want Δs —we want Δs divided by Δt . We divide the preceding equation by Δt , getting

$$\frac{\Delta s}{\Delta t} = 3At^2 + B + 3At(\Delta t) + A(\Delta t)^2.$$

Table 8-3. A Short Table of Derivatives

s, u, v, w are arbitrary functions of t ; a, b, c , and n are arbitrary constants

Function	Derivative
$s = t^n$	$\frac{ds}{dt} = nt^{n-1}$
$s = cu$	$\frac{ds}{dt} = c \frac{du}{dt}$
$s = u + v + w + \cdots$	$\frac{ds}{dt} = \frac{du}{dt} + \frac{dv}{dt} + \frac{dw}{dt} + \cdots$
$s = c$	$\frac{ds}{dt} = 0$
$s = u^a v^b w^c \cdots$	$\frac{ds}{dt} = s \left(\frac{a}{u} \frac{du}{dt} + \frac{b}{v} \frac{dv}{dt} + \frac{c}{w} \frac{dw}{dt} + \cdots \right)$

As Δt goes toward 0 the limit of $\Delta s/\Delta t$ is ds/dt and is equal to

$$\frac{ds}{dt} = 3At^2 + B.$$

This is the fundamental process of calculus, differentiating functions. The process is even more simple than it appears. Observe that when these expansions contain any term with a square or a cube or any higher power of Δt , such terms may be dropped at once, since they will go to 0 when the limit is taken. After a little practice the process gets easier because one knows what to leave out. There are many rules or formulas for differentiating various types of functions. These can be memorized, or can be found in tables. A short list is found in Table 8-3.

8-4 Distance as an integral

Now we have to discuss the inverse problem. Suppose that instead of a table of distances, we have a table of speeds at different times, starting from zero. For the falling ball, such speeds and times are shown in Table 8-4. A similar table could be constructed for the velocity of the car, by recording the speedometer reading every minute or half-minute. If we know how fast the car is going at any time, can we determine how far it goes? This problem is just the inverse of the one solved above; we are given the velocity and asked to find the distance. How can we find the distance if we know the speed? If the speed of the car is not constant, and the lady goes sixty miles an hour for a moment, then slows down, speeds up, and so on, how can we determine how far she has gone? That is easy. We use the same idea, and express the distance in terms of infinitesimals. Let us say, "In the first second her speed was such and such, and from the formula $\Delta s = v \Delta t$ we can calculate how far the car went the first second at that speed." Now in the next second her speed is nearly the same, but slightly different; we can calculate how far she went in the next second by taking the new speed times the time. We proceed similarly for each second, to the end of the run. We now have a number of little distances, and the total distance will be the sum of all these little pieces. That is, the distance will be the sum of the velocities times the times, or $s = \sum v \Delta t$, where the Greek letter \sum (sigma) is used to denote addition. To be more precise, it is the sum of the velocity at a certain time, let us say the i -th time, multiplied by Δt .

$$s = \sum_i v(t_i) \Delta t. \quad (8.6)$$

The rule for the times is that $t_{i+1} = t_i + \Delta t$. However, the distance we obtain by this method will not be correct, because the velocity changes during the time interval Δt . If we take the times short enough, the sum is precise, so we take them smaller and smaller until we obtain the desired accuracy. The true s is

$$s = \lim_{\Delta t \rightarrow 0} \sum_i v(t_i) \Delta t. \quad (8.7)$$

The mathematicians have invented a symbol for this limit, analogous to the symbol for the differential. The Δ turns into a d to remind us that the time is as small as it can be; the velocity is then called v at the time t , and the addition is written as a sum with a great "s," \int (from the Latin *summa*), which has become distorted and is now unfortunately just called an integral sign. Thus we write

$$s = \int v(t) dt. \quad (8.8)$$

This process of adding all these terms together is called integration, and it is the opposite process to differentiation. The derivative of this integral is v , so one operator (d) undoes the other (\int). One can get formulas for integrals by taking the formulas for derivatives and running them backwards, because they are related to each other inversely. Thus one can work out his own table of integrals by differentiating all sorts of functions. For every formula with a differential, we get an integral formula if we turn it around.

Table 8-4
Velocity of a Falling Ball

t (sec)	v (ft/sec)
0	0
1	32
2	64
3	96
4	128

Every function can be differentiated analytically, i.e., the process can be carried out algebraically, and leads to a definite function. But it is not possible in a simple manner to write an analytical value for any integral at will. You can calculate it, for instance, by doing the above sum, and then doing it again with a finer interval Δt and again with a finer interval until you have it nearly right. In general, given some particular function, it is not possible to find, analytically, what the integral is. One may always try to find a function which, when differentiated, gives some desired function; but one may not find it, and it may not exist, in the sense of being expressible in terms of functions that have already been given names.

8-5 Acceleration

The next step in developing the equations of motion is to introduce another idea which goes beyond the concept of velocity to that of *change* of velocity, and we now ask, "How does the velocity *change*?" In previous chapters we have discussed cases in which forces produce changes in velocity. You may have heard with great excitement about some car that can get from rest to 60 miles an hour in ten seconds flat. From such a performance we can see how fast the speed changes, but only on the average. What we shall now discuss is the next level of complexity, which is how fast the velocity is changing. In other words, by how many feet per second does the velocity change in a second, that is, how many feet per second, per second? We previously derived the formula for the velocity of a falling body as $v = 32t$, which is charted in Table 8-4, and now we want to find out how much the velocity changes per second; this quantity is called the acceleration.

Acceleration is defined as the time rate of change of velocity. From the preceding discussion we know enough already to write the acceleration as the derivative dv/dt , in the same way that the velocity is the derivative of the distance. If we now differentiate the formula $v = 32t$ we obtain, for a falling body,

$$a = \frac{dv}{dt} = 32. \quad (8.9)$$

[To differentiate the term $32t$ we can utilize the result obtained in a previous problem, where we found that the derivative of Bt is simply B (a constant). So by letting $B = 32$, we have at once that the derivative of $32t$ is 32.] This means that the velocity of a falling body is changing by 32 feet per second, per second always. We also see from Table 8-4 that the velocity increases by 32 ft/sec in each second. This is a very simple case, for accelerations are usually not constant. The reason the acceleration is constant here is that the force on the falling body is constant, and Newton's law says that the acceleration is proportional to the force.

As a further example, let us find the acceleration in the problem we have already solved for the velocity. Starting with

$$s = At^3 + Bt + C$$

we obtained, for $v = ds/dt$,

$$v = 3At^2 + B.$$

Since acceleration is the derivative of the velocity with respect to the time, we need to differentiate the last expression above. Recall the rule that the derivative of the two terms on the right equals the sum of the derivatives of the individual terms. To differentiate the first of these terms, instead of going through the fundamental process again we note that we have already differentiated a quadratic term when we differentiated $16t^2$, and the effect was to double the numerical coefficient and change the t^2 to t ; let us assume that the same thing will happen this time, and you can check the result yourself. The derivative of $3At^2$ will then be $6At$. Next we differentiate B , a constant term; but by a rule stated previously, the derivative of B is zero; hence this term contributes nothing to the acceleration. The final result, therefore, is $a = dv/dt = 6At$.

For reference, we state two very useful formulas, which can be obtained by integration. If a body starts from rest and moves with a constant acceleration, g , its velocity v at any time t is given by

$$v = gt.$$

The distance it covers in the same time is

$$s = \frac{1}{2}gt^2.$$

Various mathematical notations are used in writing derivatives. Since velocity is ds/dt and acceleration is the time derivative of the velocity, we can also write

$$a = \frac{d}{dt} \left(\frac{ds}{dt} \right) = \frac{d^2s}{dt^2}, \quad (8.10)$$

which are common ways of writing a second derivative.

We have another law that the velocity is equal to the integral of the acceleration. This is just the opposite of $a = dv/dt$; we have already seen that distance is the integral of the velocity, so distance can be found by twice integrating the acceleration.

In the foregoing discussion the motion was in only one dimension, and space permits only a brief discussion of motion in three dimensions. Consider a particle P which moves in three dimensions in any manner whatsoever. At the beginning of this chapter, we opened our discussion of the one-dimensional case of a moving car by observing the distance of the car from its starting point at various times. We then discussed velocity in terms of changes of these distances with time, and acceleration in terms of changes in velocity. We can treat three-dimensional motion analogously. It will be simpler to illustrate the motion on a two-dimensional diagram, and then extend the ideas to three dimensions. We establish a pair of axes at right angles to each other, and determine the position of the particle at any moment by measuring how far it is from each of the two axes. Thus each position is given in terms of an x -distance and a y -distance, and the motion can be described by constructing a table in which both these distances are given as functions of time. (Extension of this process to three dimensions requires only another axis, at right angles to the first two, and measuring a third distance, the z -distance. The distances are now measured from coordinate *planes* instead of lines.) Having constructed a table with x - and y -distances, how can we determine the velocity? We first find the components of velocity in each direction. The horizontal part of the velocity, or x -component, is the derivative of the x -distance with respect to the time, or

$$v_x = dx/dt. \quad (8.11)$$

Similarly, the vertical part of the velocity, or y -component, is

$$v_y = dy/dt. \quad (8.12)$$

In the third dimension,

$$v_z = dz/dt. \quad (8.13)$$

Now, given the components of velocity, how can we find the velocity along the actual path of motion? In the two-dimensional case, consider two successive positions of the particle, separated by a short distance Δs and a short time interval $t_2 - t_1 = \Delta t$. In the time Δt the particle moves horizontally a distance $\Delta x \sim v_x \Delta t$, and vertically a distance $\Delta y \sim v_y \Delta t$. (The symbol “ \sim ” is read “is approximately.”) The actual distance moved is approximately

$$\Delta s \sim \sqrt{(\Delta x)^2 + (\Delta y)^2}, \quad (8.14)$$

as shown in Fig. 8-3. The approximate velocity during this interval can be obtained by dividing by Δt and by letting Δt go to 0, as at the beginning of the chapter.

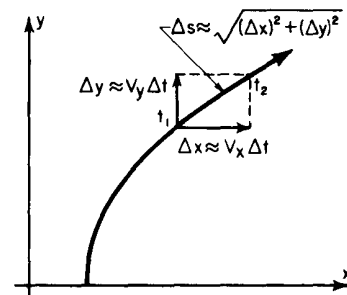


Fig. 8-3. Description of the motion of a body in two dimensions and the computation of its velocity.

We then get the velocity as

$$v = \frac{ds}{dt} = \sqrt{(dx/dt)^2 + (dy/dt)^2} = \sqrt{v_x^2 + v_y^2} \quad (8.15)$$

For three dimensions the result is

$$v = \sqrt{v_x^2 + v_y^2 + v_z^2}. \quad (8.16)$$

In the same way as we defined velocities, we can define accelerations: we have an x -component of acceleration a_x , which is the derivative of v_x , the x -component of the velocity (that is, $a_x = d^2x/dt^2$, the second derivative of x with respect to t), and so on.

Let us consider one nice example of compound motion in a plane. We shall take a motion in which a ball moves horizontally with a constant velocity u , and at the same time goes vertically downward with a constant acceleration $-g$; what is the motion? We can say $dx/dt = v_x = u$. Since the velocity v_x is constant,

$$x = ut, \quad (8.17)$$

and since the downward acceleration $-g$ is constant, the distance y the object falls can be written as

$$y = -\frac{1}{2}gt^2. \quad (8.18)$$

What is the curve of its path, i.e., what is the relation between y and x ? We can eliminate t from Eq. (8.18), since $t = x/u$. When we make this substitution we find that

$$y = -\frac{g}{2u^2} x^2. \quad (8.19)$$

This relation between y and x may be considered as the equation of the path of the moving ball. When this equation is plotted we obtain a curve that is called a parabola; any freely falling body that is shot out in any direction will travel in a parabola, as shown in Fig. 8-4.

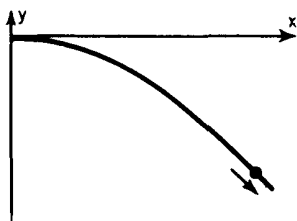


Fig. 8-4. The parabola described by a falling body with an initial horizontal velocity.

Newton's Laws of Dynamics

9-1 Momentum and force

The discovery of the laws of dynamics, or the laws of motion, was a dramatic moment in the history of science. Before Newton's time, the motions of things like the planets were a mystery, but after Newton there was complete understanding. Even the slight deviations from Kepler's laws, due to the perturbations of the planets, were computable. The motions of pendulums, oscillators with springs and weights in them, and so on, could all be analyzed completely after Newton's laws were enunciated. So it is with this chapter: before this chapter we could not calculate how a mass on a spring would move; much less could we calculate the perturbations on the planet Uranus due to Jupiter and Saturn. After this chapter we *will* be able to compute not only the motion of the oscillating mass, but also the perturbations on the planet Uranus produced by Jupiter and Saturn!

Galileo made a great advance in the understanding of motion when he discovered the *principle of inertia*: if an object is left alone, is not disturbed, it continues to move with a constant velocity in a straight line if it was originally moving, or it continues to stand still if it was just standing still. Of course this never appears to be the case in nature, for if we slide a block across a table it stops, but that is because it is *not* left to itself—it is rubbing against the table. It required a certain imagination to find the right rule, and that imagination was supplied by Galileo.

Of course, the next thing which is needed is a rule for finding how an object *changes* its speed if something *is* affecting it. That is the contribution of Newton. Newton wrote down three laws: The First Law was a mere restatement of the Galilean principle of inertia just described. The Second Law gave a specific way of determining how the velocity changes under different influences called *forces*. The Third Law describes the forces to some extent, and we shall discuss that at another time. Here we shall discuss only the Second Law, which asserts that the motion of an object is changed by forces in this way: *the time-rate-of-change of a quantity called momentum is proportional to the force*. We shall state this mathematically shortly, but let us first explain the idea.

Momentum is not the same as *velocity*. A lot of words are used in physics, and they all have precise meanings in physics, although they may not have such precise meanings in everyday language. Momentum is an example, and we must define it precisely. If we exert a certain push with our arms on an object that is light, it moves easily; if we push just as hard on another object that is much heavier in the usual sense, then it moves much less rapidly. Actually, we must change the words from "light" and "heavy" to *less massive* and *more massive*, because there is a difference to be understood between the *weight* of an object and its *inertia*. (How hard it is to get it going is one thing, and how much it weighs is something else.) Weight and inertia are *proportional*, and on the earth's surface are often taken to be numerically equal, which causes a certain confusion to the student. On Mars, weights would be different but the amount of force needed to overcome inertia would be the same.

We use the term *mass* as a quantitative measure of inertia, and we may measure mass, for example, by swinging an object in a circle at a certain speed and measuring how much force we need to keep it in the circle. In this way we find a certain quantity of mass for every object. Now the *momentum* of an object is a product of two parts: its *mass* and its *velocity*. Thus Newton's Second Law may

9-1 Momentum and force

9-2 Speed and velocity

9-3 Components of velocity, acceleration, and force

9-4 What is the force?

9-5 Meaning of the dynamical equations

9-6 Numerical solution of the equations

9-7 Planetary motions

be written mathematically this way:

$$F = \frac{d}{dt} (mv). \quad (9.1)$$

Now there are several points to be considered. In writing down any law such as this, we use many intuitive ideas, implications, and assumptions which are at first combined approximately into our "law." Later we may have to come back and study in greater detail exactly what each term means, but if we try to do this too soon we shall get confused. Thus at the beginning we take several things for granted. First, that the mass of an object is *constant*; it isn't really, but we shall start out with the Newtonian approximation that mass is constant, the same all the time, and that, further, when we put two objects together, their masses *add*. These ideas were of course implied by Newton when he wrote his equation, for otherwise it is meaningless. For example, suppose the mass varied inversely as the velocity; then the momentum would *never change* in any circumstance, so the law means nothing unless you know how the mass changes with velocity. At first we say, *it does not change*.

Then there are some implications concerning force. As a rough approximation we think of force as a kind of push or pull that we make with our muscles, but we can define it more accurately now that we have this law of motion. The most important thing to realize is that this relationship involves not only changes in the *magnitude* of the momentum or of the velocity but also in their *direction*. If the mass is constant, then Eq. (9.1) can also be written as

$$F = m \frac{dv}{dt} = ma. \quad (9.2)$$

The acceleration a is the rate of change of the velocity, and Newton's Second Law says more than that the effect of a given force varies inversely as the mass; it says also that the *direction* of the change in the velocity and the *direction* of the force are the same. Thus we must understand that a change in a velocity, or an acceleration, has a wider meaning than in common language: The velocity of a moving object can change by its speeding up, slowing down (when it slows down, we say it accelerates with a negative acceleration), or changing its direction of motion. An acceleration at right angles to the velocity was discussed in Chapter 7. There we saw that an object moving in a circle of radius R with a certain speed v along the circle falls away from a straightline path by a distance equal to $\frac{1}{2}(v^2/R)t^2$ if t is very small. Thus the formula for acceleration at right angles to the motion is

$$a = v^2/R, \quad (9.3)$$

and a force at right angles to the velocity will cause an object to move in a curved path whose radius of curvature can be found by dividing the force by the mass to get the acceleration, and then using (9.3).

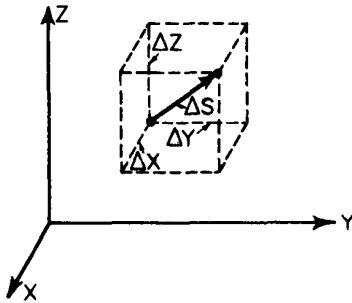


Fig. 9-1. A small displacement of an object.

9-2 Speed and velocity

In order to make our language more precise, we shall make one further definition in our use of the words *speed* and *velocity*. Ordinarily we think of speed and velocity as being the same, and in ordinary language they are the same. But in physics we have taken advantage of the fact that there *are* two words and have chosen to use them to distinguish two ideas. We carefully distinguish velocity, which has both magnitude and direction, from speed, which we choose to mean the magnitude of the velocity, but which does not include the direction. We can formulate this more precisely by describing how the x -, y -, and z -coordinates of an object change with time. Suppose, for example, that at a certain instant an object is moving as shown in Fig. 9-1. In a given small interval of time Δt it will move a certain distance Δx in the x -direction, Δy in the y -direction, and Δz in the z -direction. The total effect of these three coordinate changes is a displacement Δs along the diagonal of a parallelepiped whose sides are Δx , Δy , and Δz . In terms

of the velocity, the displacement Δx is the x -component of the velocity times Δt , and similarly for Δy and Δz :

$$\Delta x = v_x \Delta t, \quad \Delta y = v_y \Delta t, \quad \Delta z = v_z \Delta t. \quad (9.4)$$

9-3 Components of velocity, acceleration, and force

In Eq. (9.4) we have resolved the velocity into components by telling how fast the object is moving in the x -direction, the y -direction, and the z -direction. The velocity is completely specified, both as to magnitude and direction, if we give the numerical values of its three rectangular components:

$$v_x = dx/dt, \quad v_y = dy/dt, \quad v_z = dz/dt. \quad (9.5)$$

On the other hand, the speed of the object is

$$ds/dt = |v| = \sqrt{v_x^2 + v_y^2 + v_z^2}. \quad (9.6)$$

Next, suppose that, because of the action of a force, the velocity changes to some other direction and a different magnitude, as shown in Fig. 9-2. We can analyze this apparently complex situation rather simply if we evaluate the changes in the x -, y -, and z -components of velocity. The change in the component of the velocity in the x -direction in a time Δt is $\Delta v_x = a_x \Delta t$, where a_x is what we call the x -component of the acceleration. Similarly, we see that $\Delta v_y = a_y \Delta t$ and $\Delta v_z = a_z \Delta t$. In these terms, we see that Newton's Second Law, in saying that the force is in the same direction as the acceleration, is really three laws, in the sense that the component of the force in the x -, y -, or z -direction is equal to the mass times the rate of change of the corresponding component of velocity:

$$\begin{aligned} F_x &= m(dv_x/dt) = m(d^2x/dt^2) = ma_x, \\ F_y &= m(dv_y/dt) = m(d^2y/dt^2) = ma_y, \\ F_z &= m(dv_z/dt) = m(d^2z/dt^2) = ma_z. \end{aligned} \quad (9.7)$$

Just as the velocity and acceleration have been resolved into components by projecting a line segment representing the quantity and its direction onto three coordinate axes, so, in the same way, a force in a given direction is represented by certain components in the x -, y -, and z -directions:

$$\begin{aligned} F_x &= F \cos(x, F), \\ F_y &= F \cos(y, F), \\ F_z &= F \cos(z, F), \end{aligned} \quad (9.8)$$

where F is the magnitude of the force and (x, F) represents the angle between the x -axis and the direction of F , etc.

Newton's Second Law is given in complete form in Eq. (9.7). If we know the forces on an object and resolve them into x -, y -, and z -components, then we can find the motion of the object from these equations. Let us consider a simple example. Suppose there are no forces in the y - and z -directions, the only force being in the x -direction, say vertically. Equation (9.7) tells us that there would be changes in the velocity in the vertical direction, but no changes in the horizontal direction. This was demonstrated with a special apparatus in Chapter 7 (see Fig. 7-3). A falling body moves horizontally without any change in horizontal motion, while it moves vertically the same way as it would move if the horizontal motion were zero. In other words, motions in the x -, y -, and z -directions are independent if the forces are not connected.

9-4 What is the force?

In order to use Newton's laws, we have to have some formula for the force; these laws say *pay attention to the forces*. If an object is accelerating, some agency is at work; find it. Our program for the future of dynamics must be to *find the*

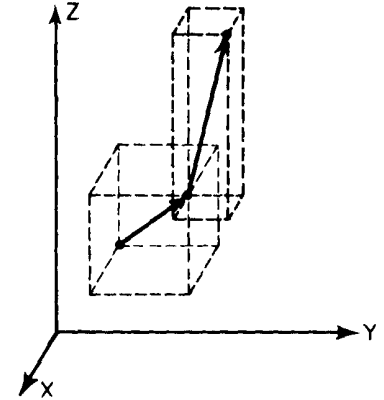


Fig. 9-2. A change in velocity in which both the magnitude and direction change.

laws for the force. Newton himself went on to give some examples. In the case of gravity he gave a specific formula for the force. In the case of other forces he gave some part of the information in his Third Law, which we will study in the next chapter, having to do with the equality of action and reaction.

Extending our previous example, what are the forces on objects near the earth's surface? Near the earth's surface, the force in the vertical direction due to gravity is proportional to the mass of the object and is nearly independent of height for heights small compared with the earth's radius R : $F = GmM/R^2 = mg$, where $g = GM/R^2$ is called the *acceleration of gravity*. Thus the law of gravity tells us that weight is proportional to mass; the force is in the vertical direction and is the mass times g . Again we find that the motion in the horizontal direction is at constant velocity. The interesting motion is in the vertical direction, and Newton's Second Law tells us

$$mg = m(d^2x/dt^2). \quad (9.9)$$

Cancelling the m 's, we find that the acceleration in the x -direction is constant and equal to g . This is of course the well known law of free fall under gravity, which leads to the equations

$$\begin{aligned} v_x &= v_0 + gt, \\ x &= x_0 + v_0t + \frac{1}{2}gt^2. \end{aligned} \quad (9.10)$$

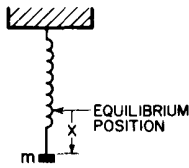


Fig. 9-3. A mass on a spring.

As another example, let us suppose that we have been able to build a gadget (Fig. 9-3) which applies a force proportional to the distance and directed oppositely—a spring. If we forget about gravity, which is of course balanced out by the initial stretch of the spring, and talk only about *excess* forces, we see that if we pull the mass down, the spring pulls up, while if we push it up the spring pulls down. This machine has been designed carefully so that the force is greater, the more we pull it up, in exact proportion to the displacement from the balanced condition, and the force upward is similarly proportional to how far we pull down. If we watch the dynamics of this machine, we see a rather beautiful motion—up, down, up, down, . . . The question is, will Newton's equations correctly describe this motion? Let us see whether we can exactly calculate how it moves with this periodic oscillation, by applying Newton's law (9.7). In the present instance, the equation is

$$-kx = m(dv_x/dt). \quad (9.11)$$

Here we have a situation where the velocity in the x -direction changes at a rate proportional to x . Nothing will be gained by retaining numerous constants, so we shall imagine either that the scale of time has changed or that there is an accident in the units, so that we happen to have $k/m = 1$. Thus we shall try to solve the equation

$$dv_x/dt = -x. \quad (9.12)$$

To proceed, we must know what v_x is, but of course we know that the velocity is the rate of change of the position.

9-5 Meaning of the dynamical equations

Now let us try to analyze just what Eq. (9.12) means. Suppose that at a given time t the object has a certain velocity v_x and position x . What is the velocity and what is the position at a slightly later time $t + \epsilon$? If we can answer this question our problem is solved, for then we can start with the given condition and compute how it changes for the first instant, the next instant, the next instant, and so on, and in this way we gradually evolve the motion. To be specific, let us suppose that at the time $t = 0$ we are given that $x = 1$ and $v_x = 0$. Why does the object move at all? Because there is a *force* on it when it is at any position except $x = 0$. If $x > 0$, that force is upward. Therefore the velocity which is zero starts to change, because of the law of motion. Once it starts to build up some velocity the object starts to move up, and so on. Now at any time t , if ϵ is very small,

we may express the position at time $t + \epsilon$ in terms of the position at time t and the velocity at time t to a very good approximation as

$$x(t + \epsilon) = x(t) + \epsilon v_x(t). \quad (9.13)$$

The smaller the ϵ , the more accurate this expression is, but it is still usefully accurate even if ϵ is not vanishingly small. Now what about the velocity? In order to get the velocity later, the velocity at the time $t + \epsilon$, we need to know how the velocity changes, the *acceleration*. And how are we going to find the acceleration? That is where the law of dynamics comes in. The law of dynamics tells us what the acceleration is. It says the acceleration is $-x$.

$$v_x(t + \epsilon) = v_x(t) + \epsilon a_x(t) \quad (9.14)$$

$$= v_x(t) - \epsilon x(t). \quad (9.15)$$

Equation (9.14) is merely kinematics; it says that a velocity changes because of the presence of acceleration. But Eq. (9.15) is *dynamics*, because it relates the acceleration to the force; it says that at this particular time for this particular problem, you can replace the acceleration by $-x(t)$. Therefore, if we know both the x and v at a given time, we know the acceleration, which tells us the new velocity, and we know the new position—this is how the machinery works. The velocity changes a little bit because of the force, and the position changes a little bit because of the velocity.

9-6 Numerical solution of the equations

Now let us really solve the problem. Suppose that we take $\epsilon = 0.100$ sec. After we do all the work if we find that this is not small enough we may have to go back and do it again with $\epsilon = 0.010$ sec. Starting with our initial value $x(0) = 1.00$, what is $x(0.1)$? It is the old position $x(0)$ plus the velocity (which is zero) times 0.10 sec. Thus $x(0.1)$ is still 1.00 because it has not yet started to move. But the new velocity at 0.10 sec will be the old velocity $v(0) = 0$ plus ϵ times the acceleration. The acceleration is $-x(0) = -1.00$. Thus

$$v(0.1) = 0.00 - 0.10 \times 1.00 = -0.10.$$

Now at 0.20 sec

$$\begin{aligned} x(0.2) &= x(0.1) + \epsilon v(0.1) \\ &= 1.00 - 0.10 \times 0.10 = 0.99 \end{aligned}$$

and

$$\begin{aligned} v(0.2) &= v(0.1) + \epsilon a(0.1) \\ &= -0.10 - 0.10 \times 1.00 = -0.20. \end{aligned}$$

And so, on and on and on, we can calculate the rest of the motion, and that is just what we shall do. However, for practical purposes there are some little tricks by which we can increase the accuracy. If we continued this calculation as we have started it, we would find the motion only rather crudely because $\epsilon = 0.100$ sec is rather crude, and we would have to go to a very small interval, say $\epsilon = 0.01$. Then to go through a reasonable total time interval would take a lot of cycles of computation. So we shall organize the work in a way that will increase the precision of our calculations, using the same coarse interval $\epsilon = 0.10$ sec. This can be done if we make a subtle improvement in the technique of the analysis.

Notice that the new position is the old position plus the time interval ϵ times the velocity. But the velocity *when*? The velocity at the beginning of the time interval is one velocity and the velocity at the end of the time interval is another velocity. Our improvement is to use the velocity *halfway between*. If we know the speed now, but the speed is changing, then we are not going to get the right answer by going at the same speed as now. We should use some speed between the “now” speed and the “then” speed at the end of the interval. The same considerations also apply to the velocity: to compute the velocity changes, we

Table 9-1

Solution of $dv_x/dt = -x$

Interval: $\epsilon = 0.10$ sec

t	x	v_x	a_x
0.0	1.000	0.000	-1.000
		-0.050	
0.1	0.995		-0.995
		-0.150	
0.2	0.980		-0.980
		-0.248	
0.3	0.955		-0.955
		-0.343	
0.4	0.921		-0.921
		-0.435	
0.5	0.877		-0.877
		-0.523	
0.6	0.825		-0.825
		-0.605	
0.7	0.764		-0.764
		-0.682	
0.8	0.696		-0.696
		-0.751	
0.9	0.621		-0.621
		-0.814	
1.0	0.540		-0.540
		-0.868	
1.1	0.453		-0.453
		-0.913	
1.2	0.362		-0.362
		-0.949	
1.3	0.267		-0.267
		-0.976	
1.4	0.169		-0.169
		-0.993	
1.5	0.070		-0.070
		-1.000	
1.6	-0.030		+0.030

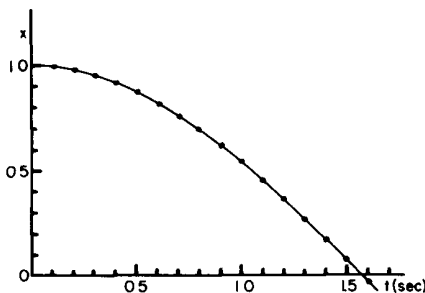


Fig. 9-4. Graph of the motion of a mass on a spring.

should use the acceleration midway between the two times at which the velocity is to be found. Thus the equations that we shall actually use will be something like this: the position later is equal to the position before plus ϵ times the velocity *at the time in the middle of the interval*. Similarly, the velocity at this halfway point is the velocity at a time ϵ before (which is in the middle of the previous interval) plus ϵ times the acceleration at the time t . That is, we use the equations

$$\begin{aligned}x(t + \epsilon) &= x(t) + \epsilon v(t + \epsilon/2), \\v(t + \epsilon/2) &= v(t - \epsilon/2) + \epsilon a(t), \\a(t) &= -x(t).\end{aligned}\tag{9.16}$$

There remains only one slight problem: what is $v(\epsilon/2)$? At the start, we are given $v(0)$, not $v(-\epsilon/2)$. To get our calculation started, we shall use a special equation, namely, $v(\epsilon/2) = v(0) + (\epsilon/2)a(0)$.

Now we are ready to carry through our calculation. For convenience, we may arrange the work in the form of a table, with columns for the time, the position, the velocity, and the acceleration, and the in-between lines for the velocity, as shown in Table 9-1. Such a table is, of course, just a convenient way of representing the numerical values obtained from the set of equations (9.16), and in fact the equations themselves need never be written. We just fill in the various spaces in the table one by one. This table now gives us a very good idea of the motion: it starts from rest, first picks up a little upward (negative) velocity and it loses some of its distance. The acceleration is then a little bit less but it is still gaining speed. But as it goes on it gains speed more and more slowly, until as it passes $x = 0$ at about $t = 1.50$ sec we can confidently predict that it will keep going, but now it will be on the other side; the position x will become negative, the acceleration therefore positive. Thus the speed decreases. It is interesting to compare these numbers with the function $x = \cos t$, which is done in Fig. 9-4. The agreement is within the three significant figure accuracy of our calculation! We shall see later that $x = \cos t$ is the exact mathematical solution of our equation of motion, but it is an impressive illustration of the power of numerical analysis that such an easy calculation should give such precise results.

9-7 Planetary motions

The above analysis is very nice for the motion of an oscillating spring, but can we analyze the motion of a planet around the sun? Let us see whether we can arrive at an approximation to an ellipse for the orbit. We shall suppose that the sun is infinitely heavy, in the sense that we shall not include its motion. Suppose a planet starts at a certain place and is moving with a certain velocity; it goes around the sun in some curve, and we shall try to analyze, by Newton's laws of motion and his law of gravitation, what the curve is. How? At a given moment it is at some position in space. If the radial distance from the sun to this position is called r , then we know that there is a force directed inward which, according to the law of gravity, is equal to a constant times the product of the sun's mass and the planet's mass divided by the square of the distance. To analyze this further we must find out what acceleration will be produced by this force. We shall need the *components* of the acceleration along two directions, which we call x and y . Thus if we specify the position of the planet at a given moment by giving x and y (we shall suppose that z is always zero because there is no force in the z -direction and, if there is no initial velocity v_z , there will be nothing to make z other than zero), the force is directed along the line joining the planet to the sun, as shown in Fig. 9-5.

From this figure we see that the horizontal component of the force is related to the complete force in the same manner as the horizontal distance x is to the complete hypotenuse r , because the two triangles are similar. Also, if x is positive, F_x is negative. That is, $F_x/|F| = -x/r$, or $F_x = -|F|x/r = -GMmx/r^3$. Now we use the dynamical law to find that this force component is equal to the mass of

the planet times the rate of change of its velocity in the x -direction. Thus we find the following laws:

$$\begin{aligned} m(dv_x/dt) &= -GMmx/r^3, \\ m(dv_y/dt) &= -GMmy/r^3, \\ r &= \sqrt{x^2 + y^2}. \end{aligned} \quad (9.17)$$

This, then, is the set of equations we must solve. Again, in order to simplify the numerical work, we shall suppose that the unit of time, or the mass of the sun, has been so adjusted (or luck is with us) that $GM \equiv 1$. For our specific example we shall suppose that the initial position of the planet is at $x = 0.500$ and $y = 0.000$, and that the velocity is all in the y -direction at the start, and is of magnitude 1.6300. Now how do we make the calculation? We again make a table with columns for the time, the x -position, the x -velocity v_x , and the x -acceleration a_x ; then, separated by a double line, three columns for position, velocity, and acceleration in the y -direction. In order to get the accelerations we are going to need Eq. (9.17); it tells us that the acceleration in the x -direction is $-x/r^3$, and the acceleration in the y -direction is $-y/r^3$, and that r is the square root of $x^2 + y^2$. Thus, given x and y , we must do a little calculating on the side, taking the square root of the sum of the squares to find r and then, to get ready to calculate the two accelerations, it is useful also to evaluate $1/r^3$. This work can be done rather easily by using a table of squares, cubes, and reciprocals: then we need only multiply x by $1/r^3$, which we do on a slide rule.

Our calculation thus proceeds by the following steps, using time intervals $\epsilon = 0.100$: Initial values at $t = 0$:

$$\begin{aligned} x(0) &= 0.500 & y(0) &= 0.000 \\ v_x(0) &= 0.000 & v_y(0) &= +1.630 \end{aligned}$$

From these we find:

$$\begin{aligned} r(0) &= 0.500 & 1/r^3(0) &= 8.000 \\ a_x &= -4.000 & a_y &= 0.000 \end{aligned}$$

Thus we may calculate the velocities $v_x(0.05)$ and $v_y(0.05)$:

$$\begin{aligned} v_x(0.05) &= 0.000 - 4.000 \times 0.050 = -0.200; \\ v_y(0.05) &= 1.630 + 0.000 \times 0.100 = 1.630. \end{aligned}$$

Now our main calculations begin:

$$\begin{aligned} x(0.1) &= 0.500 - 0.20 \times 0.1 &= 0.480 \\ y(0.1) &= 0.0 + 1.63 \times 0.1 &= 0.163 \\ r &= \sqrt{0.480^2 + 0.163^2} &= 0.507 \\ 1/r^3 &= 7.67 \\ a_x(0.1) &= 0.480 \times 7.67 &= -3.68 \\ a_y(0.1) &= -0.163 \times 7.67 &= -1.256 \\ v_x(0.15) &= -0.200 - 3.68 \times 0.1 &= -0.568 \\ v_y(0.15) &= 1.630 - 1.26 \times 0.1 &= 1.505 \\ x(0.2) &= 0.480 - 0.568 \times 0.1 &= 0.423 \\ y(0.2) &= 0.163 + 1.50 \times 0.1 &= 0.313 \\ &&&\text{etc.} \end{aligned}$$

In this way we obtain the values given in Table 9-2, and in 20 steps or so we have chased the planet halfway around the sun! In Fig. 9-6 are plotted the x - and y -coordinates given in Table 9-2. The dots represent the positions at the succession of times a tenth of a unit apart; we see that at the start the planet moves rapidly

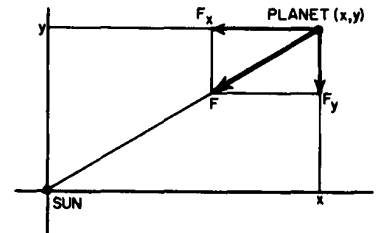


Fig. 9-5. The force of gravity on a planet.

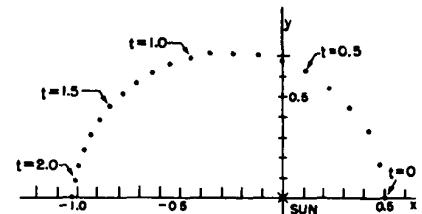


Fig. 9-6. The calculated motion of a planet around the sun.

Table 9-2

Solution of $dv_x/dt = -x/r^3$, $dv_y/dt = -y/r^3$, $r = \sqrt{x^2 + y^2}$.

Interval: $\epsilon = 0.100$

Orbit $v_y = 1.63$ $v_x = 0$ $x = 0.5$ $y = 0$ at $t = 0$

t	x	v_x	a_x	y	v_y	a_y	r	$1/r^3$
0.0	0.500		-4.00	0.000		0.00	0.500	8.000
		-0.200			1.630			
0.1	0.480		-3.68	0.163		-1.25	0.507	7.675
		-0.568			1.505			
0.2	0.423		-2.91	0.313		-2.15	0.526	6.873
		-0.859			1.290			
0.3	0.337		-1.96	0.442		-2.57	0.556	5.824
		-1.055			1.033			
0.4	0.232		-1.11	0.545		-2.62	0.592	4.81
		-1.166			0.771			
0.5	0.115		-0.453	0.622		-2.45	0.633	3.942
		-1.211			0.526			
0.6	-0.006		+0.020	0.675		-2.20	0.675	3.252
		-1.209			0.306			
0.7	-0.127		+0.344	0.706		-1.91	0.717	2.712
		-1.175			0.115			
0.8	-0.245		+0.562	0.718		-1.64	0.758	2.296
		-1.119			-0.049			
0.9	-0.357		+0.705	0.713		-1.41	0.797	1.975
		-1.048			-0.190			
1.0	-0.462		+0.796	0.694		-1.20	0.834	1.723
		-0.968			-0.310			
1.1	-0.559		+0.858	0.663		-1.02	0.867	1.535
		-0.882			-0.412			
1.2	-0.647		+0.90	0.622		-0.86	0.897	1.385
		-0.792			-0.499			
1.3	-0.726		+0.92	0.572		-0.72	0.924	1.267
		-0.700			-0.570			
1.4	-0.796		+0.93	0.515		-0.60	0.948	1.173
		-0.607			-0.630			
1.5	-0.857		+0.94	0.452		-0.50	0.969	1.099
		-0.513			-0.680			
1.6	-0.908		+0.95	0.384		-0.40	0.986	1.043
		-0.418			-0.720			
1.7	-0.950		+0.95	0.312		-0.31	1.000	1.000
		-0.323			-0.751			
1.8	-0.982		+0.95	0.237		-0.23	1.010	0.970
		-0.228			-0.773			
1.9	-1.005		+0.95	0.160		-0.15	1.018	0.948
		-0.113			-0.778			
2.0	-1.018		+0.96	0.081		-0.08	1.021	0.939
		-0.037			-0.796			
2.1	-1.022		+0.95	0.001		0.00	1.022	0.936
		+0.058			-0.796			
2.2	-1.016		+0.96	-0.079		+0.07	1.019	0.945
					-0.789			
2.3								

Crossed x-axis at 2.101 sec, \therefore period = 4.20 sec.

$v_x = 0$ at 2.086 sec.

Cross x at 1.022, \therefore semimajor axis = $\frac{1.022 + 0.500}{2} = 0.761$.

$v_y = 0.796$.

Predicted time $\pi(0.761)^{3/2} = \pi(0.663) = 2.082$.

and at the end it moves slowly, and so the shape of the curve is determined. Thus we see that we *really do* know how to calculate the motion of planets!

Now let us see how we can calculate the motion of Neptune, Jupiter, Uranus, or any other planet. If we have a great many planets, and let the sun move too, can we do the same thing? Of course we can. We calculate the force on a particular planet, let us say planet number i , which has a position x_i, y_i, z_i ($i = 1$ may represent the sun, $i = 2$ Mercury, $i = 3$ Venus, and so on). We must know the positions of all the planets. The force acting on one is due to all the other bodies which are located, let us say, at positions x_j, y_j, z_j . Therefore the equations are

$$\begin{aligned} m_i \frac{dv_{ix}}{dt} &= \sum_{j=1}^N - \frac{Gm_i m_j (x_i - x_j)}{r_{ij}^3}, \\ m_i \frac{dv_{iy}}{dt} &= \sum_{j=1}^N - \frac{Gm_i m_j (y_i - y_j)}{r_{ij}^3}, \\ m_i \frac{dv_{iz}}{dt} &= \sum_{j=1}^N - \frac{Gm_i m_j (z_i - z_j)}{r_{ij}^3}. \end{aligned} \quad (9.18)$$

Further, we define r_{ij} as the distance between the two planets i and j ; this is equal to

$$r_{ij} = \sqrt{(x_i - x_j)^2 + (y_i - y_j)^2 + (z_i - z_j)^2}. \quad (9.19)$$

Also, \sum means a sum over all values of j —all other bodies—except, of course, for $j = i$. Thus all we have to do is to make more columns, *lots* more columns. We need nine columns for the motions of Jupiter, nine for the motions of Saturn, and so on. Then when we have all initial positions and velocities we can calculate all the accelerations from Eq. (9.18) by first calculating all the distances, using Eq. (9.19). How long will it take to do it? If you do it at home, it will take a very long time! But in modern times we have machines which do arithmetic very rapidly; a very good computing machine may take 1 microsecond, that is, a millionth of a second, to do an addition. To do a multiplication takes longer, say 10 microseconds. It may be that in one cycle of calculation, depending on the problem, we may have 30 multiplications, or something like that, so one cycle will take 300 microseconds. That means that we can do 3000 cycles of computation per second. In order to get an accuracy, of, say, one part in a billion, we would need 4×10^5 cycles to correspond to one revolution of a planet around the sun. That corresponds to a computation time of 130 seconds or about two minutes. Thus it takes only two minutes to follow Jupiter around the sun, with all the perturbations of all the planets correct to one part in a billion, by this method! (It turns out that the error varies about as the square of the interval ϵ . If we make the interval a thousand times smaller, it is a million times more accurate. So, let us make the interval 10,000 times smaller.)

So, as we said, we began this chapter not knowing how to calculate even the motion of a mass on a spring. Now, armed with the tremendous power of Newton's laws, we can not only calculate such simple motions but also, given only a machine to handle the arithmetic, even the tremendously complex motions of the planets, to as high a degree of precision as we wish!

Conservation of Momentum

10-1 Newton's Third Law

On the basis of Newton's second law of motion, which gives the relation between the acceleration of any body and the force acting on it, any problem in mechanics can be solved in principle. For example, to determine the motion of a few particles, one can use the numerical method developed in the preceding chapter. But there are good reasons to make a further study of Newton's laws. First, there are quite simple cases of motion which can be analyzed not only by numerical methods, but also by direct mathematical analysis. For example, although we know that the acceleration of a falling body is 32 ft/sec^2 , and from this fact could calculate the motion by numerical methods, it is much easier and more satisfactory to analyze the motion and find the general solution, $s = s_0 + v_0 t + 16t^2$. In the same way, although we can work out the positions of a harmonic oscillator by numerical methods, it is also possible to show analytically that the general solution is a simple cosine function of t , and so it is unnecessary to go to all that arithmetical trouble when there is a simple and more accurate way to get the result. In the same manner, although the motion of one body around the sun, determined by gravitation, can be calculated point by point by the numerical methods of Chapter 9, which show the general shape of the orbit, it is nice also to get the exact shape, which analysis reveals as a perfect ellipse.

Unfortunately, there are really very few problems which can be solved exactly by analysis. In the case of the harmonic oscillator, for example, if the spring force is not proportional to the displacement, but is something more complicated, one must fall back on the numerical method. Or if there are two bodies going around the sun, so that the total number of bodies is three, then analysis cannot produce a simple formula for the motion, and in practice the problem must be done numerically. That is the famous three-body problem, which so long challenged human powers of analysis; it is very interesting how long it took people to appreciate the fact that perhaps the powers of mathematical analysis were limited and it might be necessary to use the numerical methods. Today an enormous number of problems that cannot be done analytically are solved by numerical methods, and the old three-body problem, which was supposed to be so difficult, is solved as a matter of routine in exactly the same manner that was described in the preceding chapter, namely, by doing enough arithmetic. However, there are also situations where both methods fail: the simple problems we can do by analysis, and the moderately difficult problems by numerical, arithmetical methods, but the very complicated problems we cannot do by either method. A complicated problem is, for example, the collision of two automobiles, or even the motion of the molecules of a gas. There are countless particles in a cubic millimeter of gas, and it would be ridiculous to try to make calculations with so many variables (about 10^{17} —a hundred million billion). Anything like the motion of the molecules or atoms of a gas or a block or iron, or the motion of the stars in a globular cluster, instead of just two or three planets going around the sun—such problems we cannot do directly, so we have to seek other means.

In the situations in which we cannot follow details, we need to know some general properties, that is, general theorems or principles which are consequences of Newton's laws. One of these is the principle of conservation of energy, which was discussed in Chapter 4. Another is the principle of conservation of momentum, the subject of this chapter. Another reason for studying mechanics further is that there are certain patterns of motion that are repeated in many different circum-

10-1 Newton's Third Law

10-2 Conservation of momentum

10-3 Momentum is conserved!

10-4 Momentum and energy

10-5 Relativistic momentum

stances, so it is good to study these patterns in one particular circumstance. For example, we shall study collisions; different kinds of collisions have much in common. In the flow of fluids, it does not make much difference what the fluid is, the laws of the flow are similar. Other problems that we shall study are vibrations and oscillations and, in particular, the peculiar phenomena of mechanical waves—sound, vibrations of rods, and so on.

In our discussion of Newton's laws it was explained that these laws are a kind of program that says "Pay attention to the forces," and that Newton told us only two things about the nature of forces. In the case of gravitation, he gave us the complete law of the force. In the case of the very complicated forces between atoms, he was not aware of the right laws for the forces; however, he discovered one rule, one general property of forces, which is expressed in his Third Law, and that is the total knowledge that Newton had about the nature of forces—the law of gravitation and this principle, but no other details.

This principle is that *action equals reaction*.

What is meant is something of this kind: Suppose we have two small bodies, say particles, and suppose that the first one exerts a force on the second one, pushing it with a certain force. Then, simultaneously, according to Newton's Third Law, the second particle will push on the first with an equal force, in the opposite direction; furthermore, these forces effectively act in the same line. This is the hypothesis, or law, that Newton proposed, and it seems to be quite accurate, though not exact (we shall discuss the errors later). For the moment we shall take it to be true that action equals reaction. Of course, if there is a third particle, not on the same line as the other two, the law does *not* mean that the total force on the first one is equal to the total force on the second, since the third particle, for instance, exerts its own push on each of the other two. The result is that the total effect on the first two is in some other direction, and the forces on the first two particles are, in general, neither equal nor opposite. However, the forces on each particle can be resolved into parts, there being one contribution or part due to each other interacting particle. Then each *pair* of particles has corresponding components of mutual interaction that are equal in magnitude and opposite in direction.

10-2 Conservation of momentum

Now what are the interesting consequences of the above relationship? Suppose, for simplicity, that we have just two interacting particles, possibly of different mass, and numbered 1 and 2. The forces between them are equal and opposite; what are the consequences? According to Newton's Second Law, force is the time rate of change of the momentum, so we conclude that the rate of change of momentum p_1 of particle 1 is equal to minus the rate of change of momentum p_2 of particle 2, or

$$dp_1/dt = -dp_2/dt. \quad (10.1)$$

Now if the *rate of change* is always equal and opposite, it follows that the *total change* in the momentum of particle 1 is equal and opposite to the *total change* in the momentum of particle 2; this means that if we *add* the momentum of particle 1 to the momentum of particle 2, the rate of change of the sum of these, due to the mutual forces (called internal forces) between particles, is zero; that is

$$d(p_1 + p_2)/dt = 0. \quad (10.2)$$

There is assumed to be no other force in the problem. If the rate of change of this sum is always zero, that is just another way of saying that the quantity $(p_1 + p_2)$ does not change. (This quantity is also written $m_1v_1 + m_2v_2$, and is called the *total momentum* of the two particles.) We have now obtained the result that the total momentum of the two particles does not change because of any mutual interactions between them. This statement expresses the law of conservation of

momentum in that particular example. We conclude that if there is any kind of force, no matter how complicated, between two particles, and we measure or calculate $m_1v_1 + m_2v_2$, that is, the sum of the two momenta, both before and after the forces act, the results should be equal, i.e., the total momentum is a constant.

If we extend the argument to three or more interacting particles in more complicated circumstances, it is evident that so far as internal forces are concerned, the total momentum of all the particles stays constant, since an increase in momentum of one, due to another, is exactly compensated by the decrease of the second, due to the first. That is, all the internal forces will balance out, and therefore cannot change the total momentum of the particles. Then if there are no forces from the outside (external forces), there are no forces that can change the total momentum; hence the total momentum is a constant.

It is worth describing what happens if there are forces that do *not* come from the mutual actions of the particles in question: suppose we isolate the interacting particles. If there are only mutual forces, then, as before, the total momentum of the particles does not change, no matter how complicated the forces. On the other hand, suppose there are also forces coming from the particles outside the isolated group. Any force exerted by outside bodies on inside bodies, we call an *external* force. We shall later demonstrate that the sum of all external forces equals the rate of change of the total momentum of all the particles inside, a very useful theorem.

The conservation of the total momentum of a number of interacting particles can be expressed as

$$m_1v_1 + m_2v_2 + m_3v_3 + \cdots = \text{a constant}, \quad (10.3)$$

if there are no net external forces. Here the masses and corresponding velocities of the particles are numbered 1, 2, 3, 4, . . . The general statement of Newton's Second Law for each particle,

$$f = \frac{d}{dt}(mv), \quad (10.4)$$

is true specifically for the *components* of force and momentum in any given direction; thus the *x*-component of the force on a particle is equal to the *x*-component of the rate of change of momentum of that particle, or

$$f_x = \frac{d}{dt}(mv_x), \quad (10.5)$$

and similarly for the *y*- and *z*-directions. Therefore Eq. (10.3) is really three equations, one for each direction.

In addition to the law of conservation of momentum, there is another interesting consequence of Newton's Second Law, to be proved later, but merely stated now. This principle is that the laws of physics will look the same whether we are standing still or moving with a uniform speed in a straight line. For example, a child bouncing a ball in an airplane finds that the ball bounces the same as though he were bouncing it on the ground. Even though the airplane is moving with a very high velocity, unless it changes its velocity, the laws look the same to the child as they do when the airplane is standing still. This is the so-called *relativity principle*. As we use it here we shall call it "Galilean relativity" to distinguish it from the more careful analysis made by Einstein, which we shall study later.

We have just derived the law of conservation of momentum from Newton's laws, and we could go on from here to find the special laws that describe impacts and collisions. But for the sake of variety, and also as an illustration of a kind of reasoning that can be used in physics in other circumstances where, for example, one might not know Newton's laws and might take a different approach, we shall discuss the laws of impacts and collisions from a completely different point of view. We shall base our discussion on the principle of Galilean relativity, stated above, and shall end up with the law of conservation of momentum.

We shall start by assuming that nature would look the same if we run along at a certain speed and watch it as it would if we were standing still. Before dis-

cussing collisions in which two bodies collide and stick together, or come together and bounce apart, we shall first consider two bodies that are held together by a spring or something else, and are then suddenly released and pushed by the spring or perhaps by a little explosion. Further, we shall consider motion in only one direction. First, let us suppose that the two objects are exactly the same, are nice symmetrical objects, and then we have a little explosion between them. After the explosion, one of the bodies will be moving, let us say toward the right, with a velocity v . Then it appears reasonable that the other body is moving toward the left with a velocity v , because if the objects are alike there is no reason for right or left to be preferred and so the bodies would do something that is symmetrical. This is an illustration of a kind of thinking that is very useful in many problems but would not be brought out if we just started with the formulas.

The first result from our experiment is that equal objects will have equal speed, but now suppose that we have two objects made of different materials, say copper and aluminum, and we make the two *masses* equal. We shall now suppose that if we do the experiment with two masses that are equal, even though the objects are not identical, the velocities will be equal. Someone might object: "But you know, you could do it backwards, you did not have to *suppose* that. You could *define* equal masses to mean two masses that acquire equal velocities in this experiment." We follow that suggestion and make a little explosion between the copper and a very large piece of aluminum, so heavy that the copper flies out and the aluminum hardly budes. That is too much aluminum, so we reduce the amount until there is just a very tiny piece, then when we make the explosion the aluminum goes flying away, and the *copper* hardly budes. That is not enough aluminum. Evidently there is some right amount in between; so we keep adjusting the amount until the velocities come out equal. Very well then—let us turn it around, and say that when the velocities are equal, the masses are equal. This appears to be just a definition, and it seems remarkable that we can transform physical laws into mere definitions. Nevertheless, there *are* some physical laws involved, and if we accept this definition of equal masses, we immediately find one of the laws, as follows.

Suppose we know from the foregoing experiment that two pieces of matter, A and B (of copper and aluminum), have equal masses, and we compare a third body, say a piece of gold, with the copper in the same manner as above, making sure that its mass is equal to the mass of the copper. If we now make the experiment between the aluminum and the gold, there is nothing in logic that says *these* masses must be equal; however, the *experiment* shows that they actually are. So now, by experiment, we have found a new law. A statement of this law might be: If two masses are each equal to a third mass (as determined by equal velocities in this experiment), then they are equal to each other. (This statement does *not* follow at all from a similar statement used as a postulate regarding *mathematical* quantities.) From this example we can see how quickly we start to infer things if we are careless. It is *not* just a definition to say the masses are equal when the velocities are equal, because to say the masses are equal is to imply the mathematical laws of equality, which in turn makes a prediction about an experiment.

As a second example, suppose that A and B are found to be equal by doing the experiment with one strength of explosion, which gives a certain velocity; if we then use a stronger explosion, will it be true or not true that the velocities now obtained are equal? Again, in logic there is nothing that can decide this question, but experiment shows that it *is* true. So, here is another law, which might be stated: If two bodies have equal masses, as measured by equal velocities at one velocity, they will have equal masses when measured at another velocity. From these examples we see that what appeared to be only a definition really involved some laws of physics.

In the development that follows we shall assume it is true that equal masses have equal and opposite velocities when an explosion occurs between them. We shall make another assumption in the inverse case: If two identical objects, moving in opposite directions with equal velocities, collide and stick together by some kind of glue, then which way will they be moving after the collision? This is again a

symmetrical situation, with no preference between right and left, so we assume that they stand still. We shall also suppose that any two objects of equal mass, even if the objects are made of different materials, which collide and stick together, when moving with the same velocity in opposite directions will come to rest after the collision.

10-3 Momentum is conserved!

We can verify the above assumptions experimentally: first, that if two stationary objects of equal mass are separated by an explosion they will move apart with the same speed, and second, if two objects of equal mass, coming together with the same speed, collide and stick together they will stop. This we can do by means of a marvelous invention called an air trough,* which gets rid of friction, the thing which continually bothered Galileo (Fig. 10-1). He could not do experiments by sliding things because they do not slide freely, but, by adding a magic touch, we can today get rid of friction. Our objects will slide without difficulty, on and on at a constant velocity, as advertised by Galileo. This is done by supporting the objects on air. Because air has very low friction, an object glides along with practically constant velocity when there is no applied force. First, we use two glide blocks which have been made carefully to have the same weight, or mass (their weight was measured really, but we know that this weight is proportional to the mass), and we place a small explosive cap in a closed cylinder between the two blocks (Fig. 10-2). We shall start the blocks from rest at the center point of the track and force them apart by exploding the cap with an electric spark. What should happen? If the speeds are equal when they fly apart, they should arrive at the ends of the trough at the same time. On reaching the ends they will both bounce back with practically opposite velocity, and will come together and stop at the center where they started. It is a good test; when it is actually done the result is just as we have described (Fig. 10-3).

Now the next thing we would like to figure out is what happens in a less simple situation. Suppose we have two equal masses, one moving with velocity v and the other standing still, and they collide and stick; what is going to happen? There is a mass $2m$ altogether when we are finished, drifting with an unknown velocity. What velocity? That is the problem. To find the answer, we make the assumption that if we ride along in a car, physics will look the same as if we are standing still. We start with the knowledge that two equal masses, moving in opposite directions with equal speeds v , will stop dead when they collide. Now suppose that while this happens, we are riding by in an automobile, at a velocity $-v$. Then what does it look like? Since we are riding along with one of the two masses which are coming together, that one appears to us to have zero velocity. The other mass, however, going the other way with velocity v , will appear to be coming toward us at a velocity $2v$ (Fig. 10-4). Finally, the combined masses after collision will seem to be passing by with velocity v . We therefore conclude that an object with velocity $2v$, hitting an equal one at rest, will end up with velocity v , or what is mathematically exactly the same, an object with velocity v hitting and sticking to one at rest will produce an object moving with velocity $v/2$. Note that if we multiply the mass and the velocity beforehand and add them together, $mv + 0$, we get the same answer as when we multiply the mass and the velocity of everything afterwards, $2m$ times $v/2$. So that tells us what happens when a mass of velocity v hits one standing still.

In exactly the same manner we can deduce what happens when equal objects having *any* two velocities hit each other.

Suppose we have two equal bodies with velocities v_1 and v_2 , respectively, which collide and stick together. What is their velocity v after the collision? Again we ride by in an automobile, say at velocity v_2 , so that one body appears to be at rest. The other then appears to have a velocity $v_1 - v_2$, and we have the same case that we had before. When it is all finished they will be moving at $\frac{1}{2}(v_1 - v_2)$ with respect to the car. What then is the actual speed on the ground?

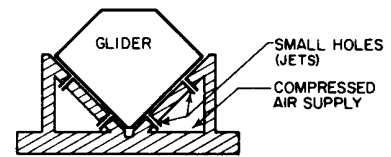


Fig. 10-1. End view of linear air trough.

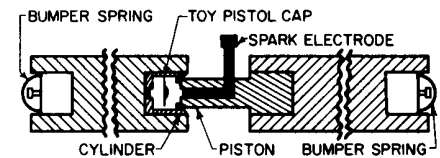


Fig. 10-2. Sectional view of gliders with explosive interaction cylinder attachment.

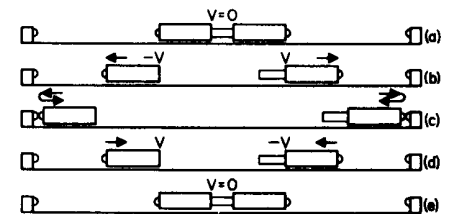


Fig. 10-3. Schematic view of action-reaction experiment with equal masses.

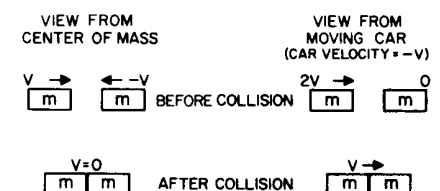


Fig. 10-4. Two views of an inelastic collision between equal masses.

* H. V. Neher and R. B. Leighton, *Amer. Jour. of Phys.* 31, 255 (1963).

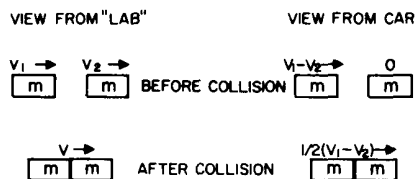


Fig. 10-5. Two views of another inelastic collision between equal masses.

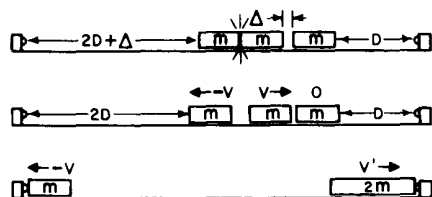


Fig. 10-6. An experiment to verify that a mass m with velocity v striking a mass m with zero velocity gives $2m$ with velocity $v/2$.

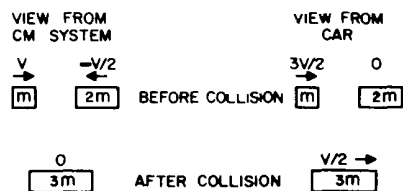


Fig. 10-7. Two views of an inelastic collision between m and $2m$.

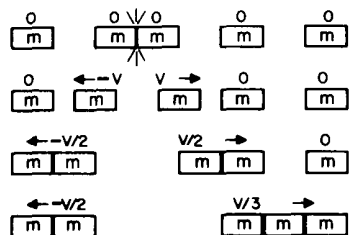


Fig. 10-8. Action and reaction between $2m$ and $3m$.

It is $v = \frac{1}{2}(v_1 - v_2) + v_2$ or $\frac{1}{2}(v_1 + v_2)$ (Fig. 10-5). Again we note that

$$mv_1 + mv_2 = 2m(v_1 + v_2)/2. \quad (10.6)$$

Thus, using this principle, we can analyze any kind of collision in which two bodies of equal mass hit each other and stick. In fact, although we have worked only in one dimension, we can find out a great deal about much more complicated collisions by imagining that we are riding by in a car in some oblique direction. The principle is the same, but the details get somewhat complicated.

In order to test experimentally whether an object moving with velocity v , colliding with an equal one at rest, forms an object moving with velocity $v/2$, we may perform the following experiment with our air-trough apparatus. We place in the trough three equally massive objects, two of which are initially joined together with our explosive cylinder device, the third being very near to but slightly separated from these and provided with a sticky bumper so that it will stick to another object which hits it. Now, a moment after the explosion, we have two objects of mass m moving with equal and opposite velocities v . A moment after that, one of these collides with the third object and makes an object of mass $2m$ moving, so we believe, with velocity $v/2$. How do we test whether it is really $v/2$? By arranging the initial positions of the masses on the trough so that the distances to the ends are not equal, but are in the ratio 2:1. Thus our first mass, which continues to move with velocity v , should cover twice as much distance in a given time as the two which are stuck together (allowing for the small distance travelled by the second object before it collided with the third). The mass m and the mass $2m$ should reach the ends at the same time, and when we try it, we find that they do (Fig. 10-6).

The next problem that we want to work out is what happens if we have two different masses. Let us take a mass m and a mass $2m$ and apply our explosive interaction. What will happen then? If, as a result of the explosion, m moves with velocity v , with what velocity does $2m$ move? The experiment we have just done may be repeated with zero separation between the second and third masses, and when we try it we get the same result, namely, the reacting masses m and $2m$ attain velocities $-v$ and $v/2$. Thus the direct reaction between m and $2m$ gives the same result as the symmetrical reaction between m and m , followed by a collision between m and a third mass m in which they stick together. Furthermore, we find that the masses m and $2m$ returning from the ends of the trough, with their velocities (nearly) exactly reversed, stop dead if they stick together.

Now the next question we may ask is this. What will happen if a mass m with velocity v , say, hits and sticks to another mass $2m$ at rest? This is very easy to answer using our principle of Galilean relativity, for we simply watch the collision which we have just described from a car moving with velocity $-v/2$ (Fig. 10-7). From the car, the velocities are

$$v'_1 = v - v(\text{car}) = v + v/2 = 3v/2$$

and

$$v'_2 = -v/2 - v(\text{car}) = -v/2 + v/2 = 0.$$

After the collision, the mass $3m$ appears to us to be moving with velocity $v/2$. Thus we have the answer, i.e., the ratio of velocities before and after collision is 3 to 1: if an object of mass m collides with a stationary object of mass $2m$, then the whole thing moves off, stuck together, with a velocity $1/3$ as much. The general rule again is that the sum of the products of the masses and the velocities stays the same: $mv + 0$ equals $3m$ times $v/3$, so we are gradually building up the theorem of the conservation of momentum, piece by piece.

Now we have one against two. Using the same arguments, we can predict the result of one against three, two against three, etc. The case of two against three, starting from rest, is shown in Fig. 10-8.

In every case we find that the mass of the first object times its velocity, plus the mass of the second object times its velocity, is equal to the total mass of the final object times its velocity. These are all examples, then, of the conservation of

momentum. Starting from simple, symmetrical cases, we have demonstrated the law for more complex cases. We could, in fact, do it for any rational mass ratio, and since every ratio is exceedingly close to a rational ratio, we can handle every ratio as precisely as we wish.

10-4 Momentum and energy

All the foregoing examples are simple cases where the bodies collide and stick together, or were initially stuck together and later separated by an explosion. However, there are situations in which the bodies do *not* cohere, as, for example, two bodies of equal mass which collide with equal speeds and then rebound. For a brief moment they are in contact and both are compressed. At the instant of maximum compression they both have zero velocity and energy is stored in the elastic bodies, as in a compressed spring. This energy is derived from the kinetic energy the bodies had before the collision, which becomes zero at the instant their velocity is zero. The loss of kinetic energy is only momentary, however. The compressed condition is analogous to the cap that releases energy in an explosion. The bodies are immediately decompressed in a kind of explosion, and fly apart again; but we already know that case—the bodies fly apart with equal speeds. However, this speed of rebound is less, in general, than the initial speed, because not all the energy is available for the explosion, depending on the material. If the material is putty no kinetic energy is recovered, but if it is something more rigid, some kinetic energy is usually regained. In the collision the rest of the kinetic energy is transformed into heat and vibrational energy—the bodies are hot and vibrating. The vibrational energy also is soon transformed into heat. It is possible to make the colliding bodies from highly elastic materials, such as steel, with carefully designed spring bumpers, so that the collision generates very little heat and vibration. In these circumstances the velocities of rebound are practically equal to the initial velocities; such a collision is called *elastic*.

That the velocities *before* and *after* an elastic collision are equal is not a matter of conservation of momentum, but a matter of conservation of *kinetic energy*. That the speeds of the bodies rebounding after a symmetrical collision are equal to *each other*, however, is a matter of conservation of momentum.

We might similarly analyze collisions between bodies of different masses, different initial velocities, and various degrees of elasticity, and determine the final velocities and the loss of kinetic energy, but we shall not go into the details of these processes.

Elastic collisions are especially interesting for systems that have no internal “gears, wheels, or parts.” Then when there is a collision there is nowhere for the energy to be impounded, because the objects that move apart are in the same condition as when they collided. Therefore, between very elementary objects, the collisions are always elastic or very nearly elastic. For instance, the collisions between atoms or molecules in a gas are said to be perfectly elastic. Although this is an excellent approximation, even such collisions are not *perfectly* elastic; otherwise one could not understand how energy in the form of light or heat radiation could come out of a gas. Once in a while, in a gas collision, a low-energy infrared ray is emitted, but this occurrence is very rare and the energy emitted is very small. So, for most purposes, collisions of molecules in gases are considered to be perfectly elastic.

As an interesting example, let us consider an *elastic* collision between two objects of *equal mass*. If they come together with the same speed, they would come apart at that same speed, by symmetry. But now look at this in another circumstance, in which one of them is moving with velocity v and the other one is at rest. What happens? We have been through this before. We watch the symmetrical collision from a car moving along with one of the objects, and we find that if a stationary body is struck elastically by another body of exactly the same mass, the moving body stops, and the one that was standing still now moves away with the same speed that the other one had; the bodies simply exchange velocities. This behavior can easily be demonstrated with a suitable impact apparatus. More

generally, if both bodies are moving, with different velocities, they simply exchange velocity at impact.

Another example of an almost elastic interaction is magnetism. If we arrange a pair of U-shaped magnets in our glide blocks, so that they repel each other, when one drifts quietly up to the other, it pushes it away and stands perfectly still, and now the other goes along, frictionlessly.

The principle of conservation of momentum is very useful, because it enables us to solve many problems without knowing the details. We did not know the details of the gas motions in the cap explosion, yet we could predict the velocities with which the bodies came apart, for example. Another interesting example is rocket propulsion. A rocket of large mass, M , ejects a small piece, of mass m , with a terrific velocity V relative to the rocket. After this the rocket, if it were originally standing still, will be moving with a small velocity, v . Using the principle of conservation of momentum, we can calculate this velocity to be

$$v = \frac{m}{M} \cdot V.$$

So long as material is being ejected, the rocket continues to pick up speed. Rocket propulsion is essentially the same as the recoil of a gun: there is no need for any air to push against.

10-5 Relativistic momentum

In modern times the law of conservation of momentum has undergone certain modifications. However, the law is still true today, the modifications being mainly in the definitions of things. In the theory of relativity it turns out that we do have conservation of momentum; the particles have mass and the momentum is still given by mv , the mass times the velocity, *but the mass changes with the velocity*, hence the momentum also changes. The mass varies with velocity according to the law

$$m = \frac{m_0}{\sqrt{1 - v^2/c^2}}, \quad (10.7)$$

where m_0 is the mass of the body at rest and c is the speed of light. It is easy to see from the formula that there is negligible difference between m and m_0 unless v is very large, and that for ordinary velocities the expression for momentum reduces to the old formula.

The components of momentum for a single particle are written as

$$p_x = \frac{m_0 v_x}{\sqrt{1 - v^2/c^2}}, \quad p_y = \frac{m_0 v_y}{\sqrt{1 - v^2/c^2}}, \quad p_z = \frac{m_0 v_z}{\sqrt{1 - v^2/c^2}}, \quad (10.8)$$

where $v^2 = v_x^2 + v_y^2 + v_z^2$. If the x -components are summed over all the interacting particles, both before and after a collision, the sums are equal; that is, momentum is conserved in the x -direction. The same holds true in any direction.

In Chapter 4 we saw that the law of conservation of energy is not valid unless we recognize that energy appears in different forms, electrical energy, mechanical energy, radiant energy, heat energy, and so on. In some of these cases, heat energy for example, the energy might be said to be "hidden." This example might suggest the question, "Are there also hidden forms of momentum—perhaps heat momentum?" The answer is that it is very hard to hide momentum for the following reasons.

The random motions of the atoms of a body furnish a measure of heat energy, if the *squares* of the velocities are summed. This sum will be a positive result, having no directional character. The heat is there, whether or not the body moves as a whole, and conservation of energy in the form of heat is not very obvious. On the other hand, if one sums the *velocities*, which have direction, and finds a result that is not zero, that means that there is a drift of the entire body in some particular direction, and such a gross momentum is readily observed. Thus there is no random internal lost momentum, because the body has net momentum only

when it moves as a whole. Therefore momentum, as a mechanical quantity, is difficult to hide. Nevertheless, momentum *can* be hidden—in the electromagnetic field, for example. This case is another effect of relativity.

One of the propositions of Newton was that interactions at a distance are instantaneous. It turns out that such is not the case; in situations involving electrical forces, for instance, if an electrical charge at one location is suddenly moved, the effects on another charge, at another place, do not appear instantaneously—there is a little delay. In those circumstances, even if the forces are equal the momentum will not check out; there will be a short time during which there will be trouble, because for a while the first charge will feel a certain reaction force, say, and will pick up some momentum, but the second charge has felt nothing and has not yet changed its momentum. It takes time for the influence to cross the intervening distance, which it does at 186,000 miles a second. In that tiny time the momentum of the particles is not conserved. Of course after the second charge has felt the effect of the first one and all is quieted down, the momentum equation will check out all right, but during that small interval momentum is not conserved. We represent this by saying that during this interval there is another kind of momentum besides that of the particle, mv , and that is momentum in the electromagnetic field. If we add the field momentum to the momentum of the particles, then momentum is conserved at any moment all the time. The fact that the electromagnetic field can possess momentum and energy makes that field very real, and so, for better understanding, the original idea that there are just the forces between particles has to be modified to the idea that a particle makes a field, and a field acts on another particle, and the field itself has such familiar properties as energy content and momentum, just as particles can have. To take another example: an electromagnetic field has waves, which we call light; it turns out that light also carries momentum with it, so when light impinges on an object it carries in a certain amount of momentum per second; this is equivalent to a force, because if the illuminated object is picking up a certain amount of momentum per second, its momentum is changing and the situation is exactly the same as if there were a force on it. Light can exert pressure by bombarding an object; this pressure is very small, but with sufficiently delicate apparatus it is measurable.

Now in quantum mechanics it turns out that momentum is a different thing—it is no longer mv . It is hard to define exactly what is meant by the velocity of a particle, but momentum still exists. In quantum mechanics the difference is that when the particles are represented as particles, the momentum is still mv , but when the particles are represented as waves, the momentum is measured by the number of waves per centimeter: the greater this number of waves, the greater the momentum. In spite of the differences, the law of conservation of momentum holds also in quantum mechanics. Even though the law $f = ma$ is false, and all the derivations of Newton were wrong for the conservation of momentum, in quantum mechanics, nevertheless, in the end, that particular law maintains itself!

Vectors

11-1 Symmetry in physics

In this chapter we introduce a subject that is technically known in physics as *symmetry in physical law*. The word “symmetry” is used here with a special meaning, and therefore needs to be defined. When is a thing symmetrical—how can we define it? When we have a picture that is symmetrical, one side is somehow the same as the other side. Professor Hermann Weyl has given this definition of symmetry: a thing is symmetrical if one can subject it to a certain operation and it appears exactly the same after the operation. For instance, if we look at a vase that is left-and-right symmetrical, then turn it 180° around the vertical axis, it looks the same. We shall adopt the definition of symmetry in Weyl’s more general form, and in that form we shall discuss symmetry of physical laws.

Suppose we build a complex machine in a certain place, with a lot of complicated interactions, and balls bouncing around with forces between them, and so on. Now suppose we build exactly the same kind of equipment at some other place, matching part by part, with the same dimensions and the same orientation, everything the same only displaced laterally by some distance. Then, if we start the two machines in the same initial circumstances, in exact correspondence, we ask: will one machine behave exactly the same as the other? Will it follow all the motions in exact parallelism? Of course the answer may well be *no*, because if we choose the wrong place for our machine it might be inside a wall and interferences from the wall would make the machine not work.

All of our ideas in physics require a certain amount of common sense in their application; they are not purely mathematical or abstract ideas. We have to understand what we mean when we say that the phenomena are the same when we move the apparatus to a new position. We mean that we move everything that we believe is relevant; if the phenomenon is not the same, we suggest that something relevant has not been moved, and we proceed to look for it. If we never find it, then we claim that the laws of physics do not have this symmetry. On the other hand, we may find it—we expect to find it—if the laws of physics do have this symmetry; looking around, we may discover, for instance, that the wall is pushing on the apparatus. The basic question is, if we define things well enough, if all the essential forces are included inside the apparatus, if all the relevant parts are moved from one place to another, will the laws be the same? Will the machinery work the same way?

It is clear that what we want to do is to move all the equipment and *essential* influences, but not *everything* in the world—planets, stars, and all—for if we do that, we have the same phenomenon again for the trivial reason that we are right back where we started. No, we cannot move *everything*. But it turns out in practice that with a certain amount of intelligence about what to move, the machinery will work. In other words, if we do not go inside a wall, if we know the origin of the outside forces, and arrange that those are moved too, then the machinery *will* work the same in one location as in another.

11-2 Translations

We shall limit our analysis to just mechanics, for which we now have sufficient knowledge. In previous chapters we have seen that the laws of mechanics can be summarized by a set of three equations for each particle:

$$m(d^2x/dt^2) = F_x, \quad m(d^2y/dt^2) = F_y, \quad m(d^2z/dt^2) = F_z. \quad (11.1)$$

11-1 Symmetry in physics

11-2 Translations

11-3 Rotations

11-4 Vectors

11-5 Vector algebra

11-6 Newton’s laws in vector notation

11-7 Scalar product of vectors

Now this means that there exists a way to *measure* x , y , and z on three perpendicular axes, and the forces along those directions, such that these laws are true. These must be measured from some origin, but *where do we put the origin?* All that Newton would tell us at first is that there *is* some place that we can measure from, perhaps the center of the universe, such that these laws are correct. But we can show immediately that we can never find the center, because if we use some other origin it would make no difference. In other words, suppose that there are two people—Joe, who has an origin in one place, and Moe, who has a parallel system whose origin is somewhere else (Fig. 11-1). Now when Joe measures the location of the point in space, he finds it at x , y , and z (we shall usually leave z out because it is too confusing to draw in a picture). Moe, on the other hand, when measuring the same point, will obtain a different x (in order to distinguish it, we will call it x'), and in principle a different y , although in our example they are numerically equal. So we have

$$x' = x - a, \quad y' = y, \quad z' = z. \quad (11.2)$$

Now in order to complete our analysis we must know what Moe would obtain for the forces. The force is supposed to act along some line, and by the force in the x -direction we mean the part of the total which is in the x -direction, which is the magnitude of the force times this cosine of its angle with the x -axis. Now we see that Moe would use exactly the same projection as Joe would use, so we have a set of equations

$$F_{x'} = F_x, \quad F_{y'} = F_y, \quad F_{z'} = F_z. \quad (11.3)$$

These would be the relationships between quantities as seen by Joe and Moe.

The question is, if Joe knows Newton's laws, and if Moe tries to write down Newton's laws, will they also be correct for him? Does it make any difference from which origin we measure the points? In other words, assuming that equations (11.1) are true, and the Eqs. (11.2) and (11.3) give the relationship of the measurements, is it or is it not true that

$$\begin{aligned} (a) \quad m(d^2x'/dt^2) &= F_{x'}, \\ (b) \quad m(d^2y'/dt^2) &= F_{y'}, \\ (c) \quad m(d^2z'/dt^2) &= F_{z'}? \end{aligned} \quad (11.4)$$

In order to test these equations we shall differentiate the formula for x' twice. First of all

$$\frac{dx'}{dt} = \frac{d}{dt}(x - a) = \frac{dx}{dt} - \frac{da}{dt}.$$

Now we shall assume that Moe's origin is fixed (not moving) relative to Joe's; therefore a is a constant and $da/dt = 0$, so we find that

$$dx'/dt = dx/dt$$

and therefore

$$d^2x'/dt^2 = d^2x/dt^2;$$

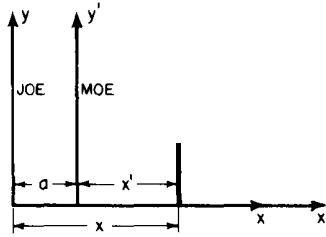
therefore we know that Eq. (11.4a) becomes

$$m(d^2x/dt^2) = F_{x'}.$$

(We also suppose that the masses measured by Joe and Moe are equal.) Thus the acceleration times the mass is the same as the other fellow's. We have also found the formula for $F_{x'}$, for, substituting from Eq. (11.1), we find that

$$F_{x'} = F_x.$$

Therefore the laws as seen by Moe appear the same; he can write Newton's laws too, with different coordinates, and they will still be right. That means that



there is no unique way to define the origin of the world, because the laws will appear the same, from whatever position they are observed.

This is also true: if there is a piece of equipment in one place with a certain kind of machinery in it, the same equipment in another place will behave in the same way. Why? Because one machine, when analyzed by Moe, has exactly the same equations as the other one, analyzed by Joe. Since the *equations* are the same, the *phenomena* appear the same. So the proof that an apparatus in a new position behaves the same as it did in the old position is the same as the proof that the equations when displaced in space reproduce themselves. Therefore we say that *the laws of physics are symmetrical for translational displacements*, symmetrical in the sense that the laws do not change when we make a translation of our coordinates. Of course it is quite obvious intuitively that this is true, but it is interesting and entertaining to discuss the mathematics of it.

11-3 Rotations

The above is the first of a series of ever more complicated propositions concerning the symmetry of a physical law. The next proposition is that it should make no difference in which *direction* we choose the axes. In other words, if we build a piece of equipment in some place and watch it operate, and nearby we build the same kind of apparatus but put it up on an angle, will it operate in the same way? Obviously it will not if it is a Grandfather clock, for example! If a pendulum clock stands upright, it works fine, but if it is tilted the pendulum falls against the side of the case and nothing happens. The theorem is then false in the case of the pendulum clock, unless we include the earth, which is pulling on the pendulum. Therefore we can make a prediction about pendulum clocks if we believe in the symmetry of physical law for rotation: something else is involved in the operation of a pendulum clock besides the machinery of the clock, something outside it that we should look for. We may also predict that pendulum clocks will not work the same way when located in different places relative to this mysterious source of asymmetry, perhaps the earth. Indeed, we know that a pendulum clock up in an artificial satellite, for example, would not tick either, because there is no effective force, and on Mars it would go at a different rate. Pendulum clocks *do* involve something more than just the machinery inside, they involve something on the outside. Once we recognize this factor, we see that we must turn the earth along with the apparatus. Of course we do not have to worry about that, it is easy to do; one simply waits a moment or two and the earth turns; then the pendulum clock ticks again in the new position the same as it did before. While we are rotating in space our angles are always changing, absolutely; this change does not seem to bother us very much, for in the new position we seem to be in the same condition as in the old. This has a certain tendency to confuse one, because it is true that in the new turned position the laws are the same as in the unturned position, but it is *not* true that *as we turn* a thing it follows the same laws as it does when we are not turning it. If we perform sufficiently delicate experiments, we can tell that the earth *is rotating*, but not that it *had rotated*. In other words, we cannot locate its angular position, but we can tell that it is changing.

Now we may discuss the effects of angular orientation upon physical laws. Let us find out whether the same game with Joe and Moe works again. This time, to avoid needless complication, we shall suppose that Joe and Moe use the same origin (we have already shown that the axes can be moved by translation to another place). Assume that Moe's axes have rotated relative to Joe's by an angle θ . The two coordinate systems are shown in Fig. 11-2, which is restricted to two dimensions. Consider any point P having coordinates (x, y) in Joe's system and (x', y') in Moe's system. We shall begin, as in the previous case, by expressing the coordinates x' and y' in terms of x, y , and θ . To do so, we first drop perpendiculars from P to all four axes and draw AB perpendicular to PQ . Inspection of the figure shows that x' can be written as the sum of two lengths along the x' -axis, and y' as the difference of two lengths along AB . All these lengths are expressed

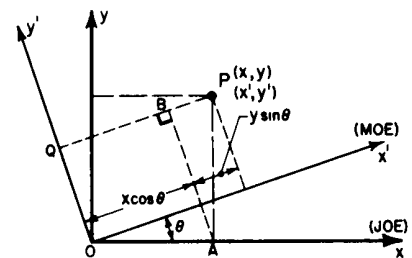


Fig. 11-2. Two coordinate systems having different angular orientations.

in terms of x , y , and θ in equations (11.5), to which we have added an equation for the third dimension.

$$\begin{aligned}x' &= x \cos \theta + y \sin \theta, \\y' &= y \cos \theta - x \sin \theta, \\z' &= z.\end{aligned}\tag{11.5}$$

The next step is to analyze the relationship of forces as seen by the two observers, following the same general method as before. Let us assume that a force F , which has already been analyzed as having components F_x and F_y (as seen by Joe), is acting on a particle of mass m , located at point P in Fig. 11-2. For simplicity, let us move both sets of axes so that the origin is at P , as shown in Fig. 11-3. Moe sees the components of F along his axes as $F_{x'}$ and $F_{y'}$. F_x has components along both the x' - and y' -axes, and F_y likewise has components along both these axes. To express $F_{x'}$ in terms of F_x and F_y , we sum these components along the x' -axis, and in a like manner we can express $F_{y'}$ in terms of F_x and F_y . The results are

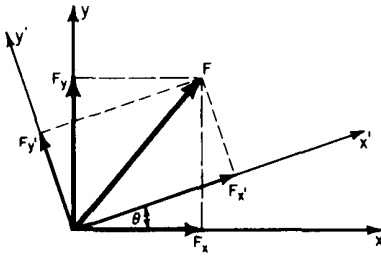


Fig. 11-3. Components of a force in the two systems.

$$\begin{aligned}F_{x'} &= F_x \cos \theta + F_y \sin \theta, \\F_{y'} &= F_y \cos \theta - F_x \sin \theta, \\F_{z'} &= F_z.\end{aligned}\tag{11.6}$$

It is interesting to note an accident of sorts, which is of extreme importance: the formulas (11.5) and (11.6), for coordinates of P and components of F , respectively, are of identical form.

As before, Newton's laws are assumed to be true in Joe's system, and are expressed by equations (11.1). The question, again, is whether Moe can apply Newton's laws—will the results be correct for his system of rotated axes? In other words, if we assume that Eqs. (11.5) and (11.6) give the relationship of the measurements, is it true or not true that

$$\begin{aligned}m(d^2x'/dt^2) &= F_{x'}, \\m(d^2y'/dt^2) &= F_{y'}, \\m(d^2z'/dt^2) &= F_{z'}?\end{aligned}\tag{11.7}$$

To test these equations, we calculate the left and right sides independently, and compare the results. To calculate the left sides, we multiply equations (11.5) by m , and differentiate twice with respect to time, assuming the angle θ to be constant. This gives

$$\begin{aligned}m(d^2x'/dt^2) &= m(d^2x/dt^2) \cos \theta + m(d^2y/dt^2) \sin \theta, \\m(d^2y'/dt^2) &= m(d^2y/dt^2) \cos \theta - m(d^2x/dt^2) \sin \theta, \\m(d^2z'/dt^2) &= m(d^2z/dt^2).\end{aligned}\tag{11.8}$$

We calculate the right sides of equations (11.7) by substituting equations (11.1) into equations (11.6). This gives

$$\begin{aligned}F_{x'} &= m(d^2x/dt^2) \cos \theta + m(d^2y/dt^2) \sin \theta, \\F_{y'} &= m(d^2y/dt^2) \cos \theta - m(d^2x/dt^2) \sin \theta, \\F_{z'} &= m(d^2z/dt^2).\end{aligned}\tag{11.9}$$

Behold! The right sides of Eqs. (11.8) and (11.9) are identical, so we conclude that if Newton's laws are correct on one set of axes, they are also valid on any other set of axes. This result, which has now been established for both translation and rotation of axes, has certain consequences: first, no one can claim his particular axes are unique, but of course they can be more *convenient* for certain particular problems. For example, it is handy to have gravity along one axis, but this is not physically necessary. Second, it means that any piece of equipment which is completely self-contained, with all the force-generating equipment completely inside the apparatus, would work the same when turned at an angle.

11-4 Vectors

Not only Newton's laws, but also the other laws of physics, so far as we know today, have the two properties which we call invariance (or symmetry) under translation of axes and rotation of axes. These properties are so important that a mathematical technique has been developed to take advantage of them in writing and using physical laws.

The foregoing analysis involved considerable tedious mathematical work. To reduce the details to a minimum in the analysis of such questions, a very powerful mathematical machinery has been devised. This system, called *vector analysis*, supplies the title of this chapter; strictly speaking, however, this is a chapter on the symmetry of physical laws. By the methods of the preceding analysis we were able to do everything required for obtaining the results that we sought, but in practice we should like to do things more easily and rapidly, so we employ the vector technique.

We began by noting some characteristics of two kinds of quantities that are important in physics. (Actually there are more than two, but let us start out with two.) One of them, like the number of potatoes in a sack, we call an ordinary quantity, or an undirected quantity, or a *scalar*. Temperature is an example of such a quantity. Other quantities that are important in physics do have direction, for instance velocity: we have to keep track of which way a body is going, not just its speed. Momentum and force also have direction, as does displacement: when someone steps from one place to another in space, we can keep track of how far he went, but if we wish also to know *where* he went, we have to specify a direction.

All quantities that have a direction, like a step in space, are called *vectors*.

A vector is three numbers. In order to represent a step in space, say from the origin to some particular point P whose location is (x, y, z) , we really need three numbers, but we are going to invent a single mathematical symbol, \mathbf{r} , which is unlike any other mathematical symbols we have so far used.* It is *not* a single number, it represents *three* numbers: x , y , and z . It means three numbers, but not really only *those* three numbers, because if we were to use a different coordinate system, the three numbers would be changed to x' , y' , and z' . However, we want to keep our mathematics simple and so we are going to use the *same mark* to represent the three numbers (x, y, z) and the three numbers (x', y', z') . That is, we use the same mark to represent the first set of three numbers for one coordinate system, but the second set of three numbers if we are using the other coordinate system. This has the advantage that when we change the coordinate system, we do not have to change the letters of our equations. If we write an equation in terms of x, y, z , and then use another system, we have to change to x', y', z' , but we shall just write \mathbf{r} , with the convention that it represents (x, y, z) if we use one set of axes, or (x', y', z') if we use another set of axes, and so on. The three numbers which describe the quantity in a given coordinate system are called the *components* of the vector in the direction of the coordinate axes of that system. That is, we use the same symbol for the three letters that correspond to the *same object, as seen from different axes*. The very fact that we can say "the same object" implies a physical intuition about the reality of a step in space, that is independent of the components in terms of which we measure it. So the symbol \mathbf{r} will represent the same thing no matter how we turn the axes.

Now suppose there is another directed physical quantity, any other quantity, which also has three numbers associated with it, like force, and these three numbers change to three other numbers by a certain mathematical rule, if we change the axes. It must be the same rule that changes (x, y, z) into (x', y', z') . In other words, any physical quantity associated with three numbers which transform as do the components of a step in space is a vector. An equation like

$$\mathbf{F} = \mathbf{r}$$

would thus be true in *any* coordinate system if it were true in one. This equation,

* In type, vectors are represented by boldface; in handwritten form an arrow is used: \vec{r} .

of course, stands for the three equations

$$F_x = x, \quad F_y = y, \quad F_z = z,$$

or, alternatively, for

$$F_{x'} = x', \quad F_{y'} = y', \quad F_{z'} = z'.$$

The fact that a physical relationship can be expressed as a vector equation assures us the relationship is unchanged by a mere rotation of the coordinate system. That is the reason why vectors are so useful in physics.

Now let us examine some of the properties of vectors. As examples of vectors we may mention velocity, momentum, force, and acceleration. For many purposes it is convenient to represent a vector quantity by an arrow that indicates the direction in which it is acting. Why can we represent force, say, by an arrow? Because it has the same mathematical transformation properties as a "step in space." We thus represent it in a diagram as if it were a step, using a scale such that one unit of force, or one newton, corresponds to a certain convenient length. Once we have done this, all forces can be represented as lengths, because an equation like

$$\mathbf{F} = k\mathbf{r},$$

where k is some constant, is a perfectly legitimate equation. Thus we can always represent forces by lines, which is very convenient, because once we have drawn the line we no longer need the axes. Of course, we can quickly calculate the three components as they change upon turning the axes, because that is just a geometric problem.

11-5 Vector algebra

Now we must describe the laws, or rules, for combining vectors in various ways. The first such combination is the *addition* of two vectors: suppose that \mathbf{a} is a vector which in some particular coordinate system has the three components (a_x, a_y, a_z) , and that \mathbf{b} is another vector which has the three components (b_x, b_y, b_z) . Now let us invent three new numbers $(a_x + b_x, a_y + b_y, a_z + b_z)$. Do these form a vector? "Well," we might say, "they are three numbers, and every three numbers form a vector." No, *not* every three numbers form a vector! In order for it to be a vector, not only must there be three numbers, but these must be associated with a coordinate system in such a way that if we turn the coordinate system, the three numbers "revolve" on each other, get "mixed up" in each other, by the precise laws we have already described. So the question is, if we now rotate the coordinate system so that (a_x, a_y, a_z) become $(a_{x'}, a_{y'}, a_{z'})$ and (b_x, b_y, b_z) become $(b_{x'}, b_{y'}, b_{z'})$, what do $(a_x + b_x, a_y + b_y, a_z + b_z)$ become? Do they become $(a_{x'} + b_{x'}, a_{y'} + b_{y'}, a_{z'} + b_{z'})$ or not? The answer is, of course, yes, because the prototype transformations of Eq. (11.5) constitute what we call a *linear* transformation. If we apply those transformations to a_x and b_x to get $a_{x'} + b_{x'}$, we find that the transformed $a_x + b_x$ is indeed the same as $a_{x'} + b_{x'}$. When \mathbf{a} and \mathbf{b} are "added together" in this sense, they will form a vector which we may call \mathbf{c} . We would write this as

$$\mathbf{c} = \mathbf{a} + \mathbf{b}.$$

Now \mathbf{c} has the interesting property

$$\mathbf{c} = \mathbf{b} + \mathbf{a},$$

as we can immediately see from its components. Thus also,

$$\mathbf{a} + (\mathbf{b} + \mathbf{c}) = (\mathbf{a} + \mathbf{b}) + \mathbf{c}.$$

We can add vectors in any order.

What is the geometric significance of $\mathbf{a} + \mathbf{b}$? Suppose that \mathbf{a} and \mathbf{b} were represented by lines on a piece of paper, what would \mathbf{c} look like? This is shown in

Fig. 11-4. We see that we can add the components of \mathbf{b} to those of \mathbf{a} most conveniently if we place the rectangle representing the components of \mathbf{b} next to that representing the components of \mathbf{a} in the manner indicated. Since \mathbf{b} just “fits” into its rectangle, as does \mathbf{a} into its rectangle, this is the same as putting the “tail” of \mathbf{b} on the “head” of \mathbf{a} , the arrow from the “tail” of \mathbf{a} to the “head” of \mathbf{b} being the vector \mathbf{c} . Of course, if we added \mathbf{a} to \mathbf{b} the other way around, we would put the “tail” of \mathbf{a} on the “head” of \mathbf{b} , and by the geometrical properties of parallelograms we would get the same result for \mathbf{c} . Note that vectors can be added in this way without reference to any coordinate axes.

Suppose we multiply a vector by a number α , what does this mean? We define it to mean a new vector whose components are αa_x , αa_y , and αa_z . We leave it as a problem for the student to prove that it is a vector.

Now let us consider vector subtraction. We may define subtraction in the same way as addition, but instead of adding, we subtract the components. Or we might define subtraction by defining a negative vector, $-\mathbf{b} = -1\mathbf{b}$, and then we would add the components. It comes to the same thing. The result is shown in Fig. 11-5. This figure shows $\mathbf{d} = \mathbf{a} - \mathbf{b} = \mathbf{a} + (-\mathbf{b})$; we also note that the difference $\mathbf{a} - \mathbf{b}$ can be found very easily from \mathbf{a} and \mathbf{b} by using the equivalent relation $\mathbf{a} = \mathbf{b} + \mathbf{d}$. Thus the difference is even easier to find than the sum: we just draw the vector from \mathbf{b} to \mathbf{a} , to get $\mathbf{a} - \mathbf{b}$!

Next we discuss velocity. Why is velocity a vector? If position is given by the three coordinates (x, y, z) , what is the velocity? The velocity is given by dx/dt , dy/dt , and dz/dt . Is that a vector, or not? We can find out by differentiating the expressions in Eq. (11.5) to find out whether dx'/dt transforms in the right way. We see that the components dx/dt and dy/dt do transform according to the same law as x and y , and therefore the time derivative is a vector. So the velocity is a vector. We can write the velocity in an interesting way as

$$\mathbf{v} = d\mathbf{r}/dt.$$

What the velocity is, and why it is a vector, can also be understood more pictorially: How far does a particle move in a short time Δt ? Answer: $\Delta \mathbf{r}$, so if a particle is “here” at one instant and “there” at another instant, then the vector difference of the positions $\Delta \mathbf{r} = \mathbf{r}_2 - \mathbf{r}_1$, which is in the direction of motion shown in Fig. 11-6, divided by the time interval $\Delta t = t_2 - t_1$, is the “average velocity” vector.

In other words, by vector velocity we mean the limit, as Δt goes to 0, of the difference between the radius vectors at the time $t + \Delta t$ and the time t , divided by Δt :

$$\mathbf{v} = \lim_{\Delta t \rightarrow 0} (\Delta \mathbf{r} / \Delta t) = d\mathbf{r} / dt. \quad (11.10)$$

Thus velocity is a vector because it is the difference of two vectors. It is also the right definition of velocity because its components are dx/dt , dy/dt , and dz/dt . In fact, we see from this argument that if we differentiate *any* vector with respect to time we produce a new vector. So we have several ways of producing new vectors: (1) multiply by a constant, (2) differentiate with respect to time, (3) add or subtract two vectors.

11-6 Newton's laws in vector notation

In order to write Newton's laws in vector form, we have to go just one step further, and define the acceleration vector. This is the time derivative of the velocity vector, and it is easy to demonstrate that its components are the second derivatives of x , y , and z with respect to t :

$$\mathbf{a} = \frac{d\mathbf{v}}{dt} = \left(\frac{d}{dt} \right) \left(\frac{d\mathbf{r}}{dt} \right) = \frac{d^2\mathbf{r}}{dt^2}, \quad (11.11)$$

$$a_x = \frac{dv_x}{dt} = \frac{d^2x}{dt^2}, \quad a_y = \frac{dv_y}{dt} = \frac{d^2y}{dt^2}, \quad a_z = \frac{dv_z}{dt} = \frac{d^2z}{dt^2}. \quad (11.12)$$

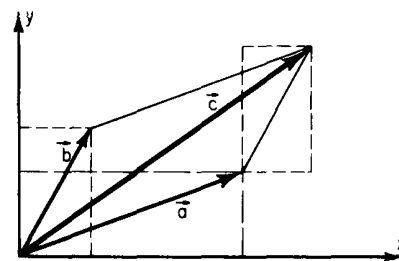


Fig. 11-4. The addition of vectors.

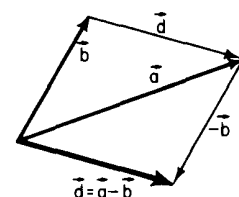


Fig. 11-5. The subtraction of vectors.

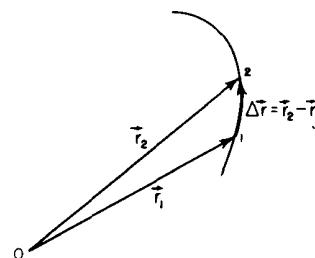


Fig. 11-6. The displacement of a particle in a short time interval $\Delta t = t_2 - t_1$.

With this definition, then, Newton's laws can be written in this way:

$$m\mathbf{a} = \mathbf{F} \quad (11.13)$$

or

$$m(d^2\mathbf{r}/dt^2) = \mathbf{F}. \quad (11.14)$$

Now the problem of proving the invariance of Newton's laws under rotation of coordinates is this: prove that \mathbf{a} is a vector; this we have just done. Prove that \mathbf{F} is a vector; we *suppose* it is. So if force is a vector, then, since we know acceleration is a vector, Eq. (11.13) will look the same in any coordinate system. Writing it in a form which does not explicitly contain x 's, y 's, and z 's has the advantage that from now on we need not write *three* laws every time we write Newton's equations or other laws of physics. We write what looks like *one* law, but really, of course, it is the three laws for any particular set of axes, because any vector equation involves the statement that *each of the components is equal*.

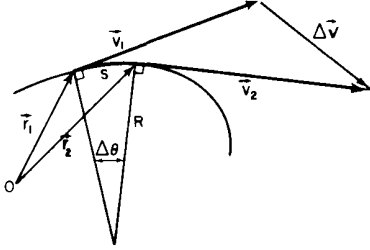


Fig. 11-7. A curved trajectory.

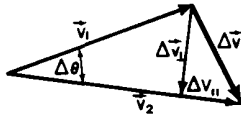


Fig. 11-8. Diagram for calculating the acceleration.

The fact that the acceleration is the rate of change of the vector velocity helps us to calculate the acceleration in some rather complicated circumstances. Suppose, for instance, that a particle is moving on some complicated curve (Fig. 11-7) and that, at a given instant t , it had a certain velocity \mathbf{v}_1 , but that when we go to another instant t_2 a little later, it has a different velocity \mathbf{v}_2 . What is the acceleration? Answer: Acceleration is the difference in the velocity divided by the small time interval, so we need the difference of the two velocities. How do we get the difference of the velocities? To subtract two vectors, we put the vector across the ends of \mathbf{v}_2 and \mathbf{v}_1 ; that is, we draw Δ as the difference of the two vectors, right? *No!* That only works when the *tails* of the vectors are in the same place! It has no meaning if we move the vector somewhere else and then draw a line across, so watch out! We have to draw a new diagram to subtract the vectors. In Fig. 11-8, \mathbf{v}_1 and \mathbf{v}_2 are both drawn parallel and equal to their counterparts in Fig. 11-7, and now we can discuss the acceleration. Of course the acceleration is simply $\Delta\mathbf{v}/\Delta t$. It is interesting to note that we can compose the velocity difference out of two parts; we can think of acceleration as having *two components*, $\Delta\mathbf{v}_{\parallel}$ in the direction tangent to the path and $\Delta\mathbf{v}_{\perp}$ at right angles to the path, as indicated in Fig. 11-8. The acceleration tangent to the path is, of course, just the change in the *length* of the vector, i.e., the change in the *speed* v :

$$a_{\parallel} = dv/dt. \quad (11.15)$$

The other component of acceleration, at right angles to the curve, is easy to calculate, using Figs. 11-7 and 11-8. In the short time Δt let the change in angle between \mathbf{v}_1 and \mathbf{v}_2 be the small angle $\Delta\theta$. If the magnitude of the velocity is called v , then of course

$$\Delta v_{\perp} = v \Delta\theta$$

and the acceleration a will be

$$a_{\perp} = v (\Delta\theta/\Delta t).$$

Now we need to know $\Delta\theta/\Delta t$, which can be found this way: If, at the given moment, the curve is approximated as a circle of a certain radius R , then in a time Δt the distance s is, of course, $v \Delta t$, where v is the speed.

$$\Delta\theta = v (\Delta t/R), \quad \text{or} \quad \Delta\theta/\Delta t = v/R.$$

Therefore, we find

$$a = v^2/R, \quad (11.16)$$

as we have seen before.

11-7 Scalar product of vectors

Now let us examine a little further the properties of vectors. It is easy to see that the *length* of a step in space would be the same in any coordinate system. That is, if a particular step \mathbf{r} is represented by x , y , z , in one coordinate system,

and by x' , y' , z' in another coordinate system, surely the distance $r = |\mathbf{r}|$ would be the same in both. Now

$$r = \sqrt{x^2 + y^2 + z^2}$$

and also

$$r' = \sqrt{x'^2 + y'^2 + z'^2}.$$

So what we wish to verify is that these two quantities are equal. It is much more convenient not to bother to take the square root, so let us talk about the square of the distance; that is, let us find out whether

$$x^2 + y^2 + z^2 = x'^2 + y'^2 + z'^2. \quad (11.17)$$

It had better be—and if we substitute Eq. (11.5) we do indeed find that it is. So we see that there are other kinds of equations which are true for any two coordinate systems.

Something new is involved. We can produce a new quantity, a function of x , y , and z , called a *scalar function*, a quantity which has no direction but which is the same in both systems. Out of a vector we can make a scalar. We have to find a general rule for that. It is clear what the rule is for the case just considered: add the squares of the components. Let us now define a new thing, which we call $\mathbf{a} \cdot \mathbf{a}$. This is not a vector, but a scalar; it is a number that is the same in all coordinate systems, and it is defined to be the sum of the squares of the three components of the vector:

$$\mathbf{a} \cdot \mathbf{a} = a_x^2 + a_y^2 + a_z^2. \quad (11.18)$$

Now you say, “But with what axes?” It does not depend on the axes, the answer is the same in *every* set of axes. So we have a new *kind* of quantity, a new *invariant* or *scalar* produced by one vector “squared.” If we now define the following quantity for any two vectors \mathbf{a} and \mathbf{b} :

$$\mathbf{a} \cdot \mathbf{b} = a_x b_x + a_y b_y + a_z b_z, \quad (11.19)$$

we find that this quantity, calculated in the primed and unprimed systems, also stays the same. To prove it we note that it is true of $\mathbf{a} \cdot \mathbf{a}$, $\mathbf{b} \cdot \mathbf{b}$, and $\mathbf{c} \cdot \mathbf{c}$, where $\mathbf{c} = \mathbf{a} + \mathbf{b}$. Therefore the sum of the squares $(a_x + b_x)^2 + (a_y + b_y)^2 + (a_z + b_z)^2$ will be invariant:

$$(a_x + b_x)^2 + (a_y + b_y)^2 + (a_z + b_z)^2 = (a_{x'} + b_{x'})^2 + (a_{y'} + b_{y'})^2 + (a_{z'} + b_{z'})^2. \quad (11.20)$$

If both sides of this equation are expanded, there will be cross products of just the type appearing in Eq. (11.19), as well as the sums of squares of the components of \mathbf{a} and \mathbf{b} . The invariance of terms of the form of Eq. (11.18) then leaves the cross product terms (11.19) invariant also.

The quantity $\mathbf{a} \cdot \mathbf{b}$ is called the *scalar product* of two vectors, \mathbf{a} and \mathbf{b} , and it has many interesting and useful properties. For instance, it is easily proved that

$$\mathbf{a} \cdot (\mathbf{b} + \mathbf{c}) = \mathbf{a} \cdot \mathbf{b} + \mathbf{a} \cdot \mathbf{c}. \quad (11.21)$$

Also, there is a simple geometrical way to calculate $\mathbf{a} \cdot \mathbf{b}$, without having to calculate the components of \mathbf{a} and \mathbf{b} : $\mathbf{a} \cdot \mathbf{b}$ is the product of the length of \mathbf{a} and the length of \mathbf{b} times the cosine of the angle between them. Why? Suppose that we choose a special coordinate system in which the x -axis lies along \mathbf{a} ; in those circumstances, the only component of \mathbf{a} that will be there is a_x , which is of course the whole length of \mathbf{a} . Thus Eq. (11.19) reduces to $\mathbf{a} \cdot \mathbf{b} = a_x b_x$ for this case, and this is the length of \mathbf{a} times the component of \mathbf{b} in the direction of \mathbf{a} , that is, $b \cos \theta$:

$$\mathbf{a} \cdot \mathbf{b} = ab \cos \theta.$$

Therefore, in that special coordinate system, we have proved that $\mathbf{a} \cdot \mathbf{b}$ is the

length of \mathbf{a} times the length of \mathbf{b} times $\cos \theta$. But *if it is true in one coordinate system, it is true in all*, because $\mathbf{a} \cdot \mathbf{b}$ is independent of the coordinate system; that is our argument.

What good is the dot product? Are there any cases in physics where we need it? Yes, we need it all the time. For instance, in Chapter 4 the kinetic energy was called $\frac{1}{2}mv^2$, but if the object is moving in space it should be the velocity squared in the x -direction, the y -direction, and the z -direction, and so the formula for kinetic energy according to vector analysis is

$$\text{K.E.} = \frac{1}{2}m(\mathbf{v} \cdot \mathbf{v}) = \frac{1}{2}m(v_x^2 + v_y^2 + v_z^2). \quad (11.22)$$

Energy does not have direction. Momentum has direction; it is a vector, and it is the mass times the velocity vector.

Another example of a dot product is the work done by a force when something is pushed from one place to the other. We have not yet defined work, but it is equivalent to the energy change, the weights lifted, when a force \mathbf{F} acts through a distance \mathbf{s} :

$$\text{Work} = \mathbf{F} \cdot \mathbf{s} \quad (11.23)$$

It is sometimes very convenient to talk about the component of a vector in a certain direction (say the vertical direction because that is the direction of gravity). For such purposes, it is useful to invent what we call a *unit vector* in the direction that we want to study. By a unit vector we mean one whose dot product with itself is equal to unity. Let us call this unit vector \mathbf{i} ; then $\mathbf{i} \cdot \mathbf{i} = 1$. Then, if we want the component of some vector in the direction of \mathbf{i} , we see that the dot product $\mathbf{a} \cdot \mathbf{i}$ will be $a \cos \theta$, i.e., the component of \mathbf{a} in the direction of \mathbf{i} . This is a nice way to get the component; in fact, it permits us to get *all* the components and to write a rather amusing formula. Suppose that in a given system of coordinates, x , y , and z , we invent three vectors: \mathbf{i} , a unit vector in the direction x ; \mathbf{j} , a unit vector in the direction y ; and \mathbf{k} , a unit vector in the direction z . Note first that $\mathbf{i} \cdot \mathbf{i} = 1$. What is $\mathbf{i} \cdot \mathbf{j}$? When two vectors are at right angles, their dot product is zero. Thus

$$\begin{aligned} \mathbf{i} \cdot \mathbf{i} &= 1 \\ \mathbf{i} \cdot \mathbf{j} &= 0 & \mathbf{j} \cdot \mathbf{j} &= 1 \\ \mathbf{i} \cdot \mathbf{k} &= 0 & \mathbf{j} \cdot \mathbf{k} &= 0 & \mathbf{k} \cdot \mathbf{k} &= 1 \end{aligned} \quad (11.24)$$

Now with these definitions, any vector whatsoever can be written this way:

$$\mathbf{a} = a_x \mathbf{i} + a_y \mathbf{j} + a_z \mathbf{k}. \quad (11.25)$$

By this means we can go from the components of a vector to the vector itself.

This discussion of vectors is by no means complete. However, rather than try to go more deeply into the subject now, we shall first learn to use in physical situations some of the ideas so far discussed. Then, when we have properly mastered this basic material, we shall find it easier to penetrate more deeply into the subject without getting too confused. We shall later find that it is useful to define another kind of product of two vectors, called the vector product, and written as $\mathbf{a} \times \mathbf{b}$. However, we shall undertake a discussion of such matters in a later chapter.

Characteristics of Force

12-1 What is a force?

Although it is interesting and worth while to study the physical laws simply because they help us to understand and to use nature, one ought to stop every once in a while and think, "What do they really mean?" The meaning of any statement is a subject that has interested and troubled philosophers from time immemorial, and the meaning of physical laws is even more interesting, because it is generally believed that these laws represent some kind of real knowledge. The meaning of knowledge is a deep problem in philosophy, and it is always important to ask, "What does it mean?"

Let us ask, "What is the meaning of the physical laws of Newton, which we write as $F = ma$? What is the meaning of force, mass, and acceleration?" Well, we can intuitively sense the meaning of mass, and we can *define* acceleration if we know the meaning of position and time. We shall not discuss those meanings, but shall concentrate on the new concept of *force*. The answer is equally simple: "If a body is accelerating, then there is a force on it." That is what Newton's laws say, so the most precise and beautiful definition of force imaginable might simply be to say that force is the mass of an object times the acceleration. Suppose we have a law which says that the conservation of momentum is valid if the sum of all the external forces is zero; then the question arises, "What does it *mean*, that the sum of all the external forces is zero?" A pleasant way to define that statement would be: "When the total momentum is a constant, then the sum of the external forces is zero." There must be something wrong with that, because it is just not saying anything new. If we have discovered a fundamental law, which asserts that the force is equal to the mass times the acceleration, and then *define* the force to be the mass times the acceleration, we have found out nothing. We could also define force to mean that a moving object with no force acting on it continues to move with constant velocity in a straight line. If we then observe an object *not* moving in a straight line with a constant velocity, we might say that there is a force on it. Now such things certainly cannot be the content of physics, because they are definitions going in a circle. The Newtonian statement above, however, seems to be a most precise definition of force, and one that appeals to the mathematician; nevertheless, it is completely useless, because no prediction whatsoever can be made from a definition. One might sit in an armchair all day long and define words at will, but to find out what happens when two balls push against each other, or when a weight is hung on a spring, is another matter altogether, because the way the bodies *behave* is something completely outside any choice of definitions.

For example, if we were to choose to say that an object left to itself keeps its position and does not move, then when we see something drifting, we could say that must be due to a "gorce"—a gorce is the rate of change of position. Now we have a wonderful new law, everything stands still except when a gorce is acting. You see, that would be analogous to the above definition of force, and it would contain no information. The real content of Newton's laws is this: that the force is supposed to have some *independent properties*, in addition to the law $F = ma$; but the *specific* independent properties that the force has were not completely described by Newton or by anybody else, and therefore the physical law $F = ma$ is an incomplete law. It implies that if we study the mass times the acceleration and call the product the force, i.e., if we study the characteristics of force as a program

12-1 What is a force?

12-2 Friction

12-3 Molecular forces

12-4 Fundamental forces. Fields

12-5 Pseudo forces

12-6 Nuclear forces

of interest, then we shall find that forces have some simplicity; the law is a good program for analyzing nature, it is a suggestion that the forces will be simple.

Now the first example of such forces was the complete law of gravitation, which was given by Newton, and in stating the law he answered the question, "What is the force?" If there were nothing but gravitation, then the combination of this law and the force law (second law of motion) would be a complete theory, but there is much more than gravitation, and we want to use Newton's laws in many different situations. Therefore in order to proceed we have to tell something about the properties of force.

For example, in dealing with force the tacit assumption is always made that the force is equal to zero unless some physical body is present, that if we find a force that is not equal to zero we also find something in the neighborhood that is a source of the force. This assumption is entirely different from the case of the "gorce" that we introduced above. One of the most important characteristics of force is that it has a material origin, and this is *not* just a definition.

Newton also gave one rule about the force: that the forces between interacting bodies are equal and opposite—action equals reaction; that rule, it turns out, is not exactly true. In fact, the law $F = ma$ is not exactly true; if it were a definition we should have to say that it is *always* exactly true; but it is not.

The student may object, "I do not like this imprecision, I should like to have everything defined exactly; in fact, it says in some books that any science is an exact subject, in which *everything* is defined." If you insist upon a precise definition of force, you will never get it! First, because Newton's Second Law is not exact, and second, because in order to understand physical laws you must understand that they are all some kind of approximation.

Any simple idea is approximate; as an illustration, consider an object, . . . what *is* an object? Philosophers are always saying, "Well, just take a chair for example." The moment they say that, you know that they do not know what they are talking about any more. What *is* a chair? Well, a chair is a certain thing over there . . . certain?, how certain? The atoms are evaporating from it from time to time—not many atoms, but a few—dirt falls on it and gets dissolved in the paint; so to define a chair precisely, to say exactly which atoms are chair, and which atoms are air, or which atoms are dirt, or which atoms are paint that belongs to the chair is impossible. So the mass of a chair can be defined only approximately. In the same way, to define the mass of a single object is impossible, because there are not any single, left-alone objects in the world—every object is a mixture of a lot of things, so we can deal with it only as a series of approximations and idealizations.

The trick is the idealizations. To an excellent approximation of perhaps one part in 10^{10} , the number of atoms in the chair does not change in a minute, and if we are not too precise we may idealize the chair as a definite thing; in the same way we shall learn about the characteristics of force, in an ideal fashion, if we are not too precise. One may be dissatisfied with the approximate view of nature that physics tries to obtain (the attempt is always to increase the accuracy of the approximation), and may prefer a mathematical definition; but mathematical definitions can never work in the real world. A mathematical definition will be good for mathematics, in which all the logic can be followed out completely, but the physical world is complex, as we have indicated in a number of examples, such as those of the ocean waves and a glass of wine. When we try to isolate pieces of it, to talk about one mass, the wine and the glass, how can we know which is which, when one dissolves in the other? The forces on a single thing already involve approximation, and if we have a system of discourse about the real world, then that system, at least for the present day, must involve approximations of some kind.

This system is quite unlike the case of mathematics, in which everything can be defined, and then we do not *know* what we are talking about. In fact, the glory of mathematics is that *we do not have to say what we are talking about*. The glory is that the laws, the arguments, and the logic are independent of what "it" is. If we have any other set of objects that obey the same system of axioms as Euclid's

geometry, then if we make new definitions and follow them out with correct logic, all the consequences will be correct, and it makes no difference what the subject was. In nature, however, when we draw a line or establish a line by using a light beam and a theodolite, as we do in surveying, are we measuring a line in the sense of Euclid? No, we are making an approximation; the cross hair has some width, but a geometrical line has no width, and so, whether Euclidean geometry can be used for surveying or not is a physical question, not a mathematical question. However, from an experimental standpoint, not a mathematical standpoint, we need to know whether the laws of Euclid apply to the kind of geometry that we use in measuring land; so we make a hypothesis that it does, and it works pretty well; but it is not precise, because our surveying lines are not really geometrical lines. Whether or not those lines of Euclid, which are really abstract, apply to the lines of experience is a question for experience; it is not a question that can be answered by sheer reason.

In the same way, we cannot just call $F = ma$ a definition, deduce everything purely mathematically, and make mechanics a mathematical theory, when mechanics is a description of nature. By establishing suitable postulates it is always possible to make a system of mathematics, just as Euclid did, but we cannot make a mathematics of the world, because sooner or later we have to find out whether the axioms are valid for the objects of nature. Thus we immediately get involved with these complicated and "dirty" objects of nature, but with approximations ever increasing in accuracy.

12-2 Friction

The foregoing considerations show that a true understanding of Newton's laws requires a discussion of forces, and it is the purpose of this chapter to introduce such a discussion, as a kind of completion of Newton's laws. We have already studied the definitions of acceleration and related ideas, but now we have to study the properties of force, and this chapter, unlike the previous chapters, will not be very precise, because forces are quite complicated.

To begin with a particular force, let us consider the drag on an airplane flying through the air. What is the law for that force? (Surely there is a law for every force, we *must* have a law!) One can hardly think that the law for that force will be simple. Try to imagine what makes a drag on an airplane flying through the air—the air rushing over the wings, the swirling in the back, the changes going on around the fuselage, and many other complications, and you see that there is not going to be a simple law. On the other hand, it is a remarkable fact that the drag force on an airplane is approximately a constant times the square of the velocity, or $F \sim cv^2$.

Now what is the status of such a law, is it analogous to $F = ma$? Not at all, because in the first place this law is an empirical thing that is obtained roughly by tests in a wind tunnel. You say, "Well $F = ma$ might be empirical too." That is not the reason that there is a difference. The difference is not that it is empirical, but that, as we understand nature, this law is the result of an enormous complexity of events and is not, fundamentally, a simple thing. If we continue to study it more and more, measuring more and more accurately, the law will continue to become *more* complicated, not *less*. In other words, as we study this law of the drag on an airplane more and more closely, we find out that it is "falser" and "falser," and the more deeply we study it, and the more accurately we measure, the more complicated the truth becomes; so in that sense we consider it not to result from a simple, fundamental process, which agrees with our original surmise. For example, if the velocity is extremely low, so low that an ordinary airplane is not flying, as when the airplane is dragged slowly through the air, then the law changes, and the drag friction depends more nearly linearly on the velocity. To take another example, the frictional drag on a ball or a bubble or anything that is moving slowly through a viscous liquid like honey, is proportional to the velocity, but for motion so fast that the fluid swirls around (honey does not but water and air do) then the drag becomes more nearly proportional to the square of the velocity ($F = cv^2$), and

if the velocity continues to increase, then even this law begins to fail. People who say, "Well the coefficient changes slightly," are dodging the issue. Second, there are other great complications: can this force on the airplane be divided or analyzed as a force on the wings, a force on the front, and so on? Indeed, this can be done, if we are concerned about the torques here and there, but then we have to get special laws for the force on the wings, and so on. It is an amazing fact that the force on a wing depends upon the other wing: in other words, if we take the airplane apart and put just one wing in the air, then the force is not the same as if the rest of the plane were there. The reason, of course, is that some of the wind that hits the front goes around to the wings and changes the force on the wings. It seems a miracle that there is such a simple, rough, empirical law that can be used in the design of airplanes, but this law is not in the same class as the *basic* laws of physics, and further study of it will only make it more and more complicated. A study of how the coefficient c depends on the shape of the front of the airplane is, to put it mildly, frustrating. There just is no simple law for determining the coefficient in terms of the shape of the airplane. In contrast, the law of gravitation is simple, and further study only indicates its greater simplicity.

We have just discussed two cases of friction, resulting from fast movement in air and slow movement in honey. There is another kind of friction, called dry friction or sliding friction, which occurs when one solid body slides on another. In this case a force is needed to maintain motion. This is called a frictional force, and its origin, also, is a very complicated matter. Both surfaces of contact are irregular, on an atomic level. There are many points of contact where the atoms seem to cling together, and then, as the sliding body is pulled along, the atoms snap apart and vibration ensues; something like that has to happen. Formerly the mechanism of this friction was thought to be very simple, that the surfaces were merely full of irregularities and the friction originated in lifting the slider over the bumps; but this cannot be, for there is no loss of energy in that process, whereas power is in fact consumed. The mechanism of power loss is that as the slider snaps over the bumps, the bumps deform and then generate waves and atomic motions and, after a while, heat, in the two bodies. Now it is very remarkable that again, empirically, this friction can be described approximately by a simple law. This law is that the force needed to overcome friction and to drag one object over another depends upon the normal force (i.e., perpendicular to the surface) between the two surfaces that are in contact. Actually, to a fairly good approximation, the frictional force is proportional to this normal force, and has a more or less constant coefficient; that is,

$$F = \mu N, \quad (12.1)$$

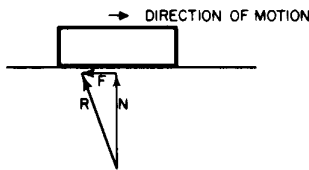


Fig. 12-1. The relation between frictional force and the normal force for sliding contact.

where μ is called the *coefficient of friction* (Fig. 12-1). Although this coefficient is not exactly constant, the formula is a good empirical rule for judging approximately the amount of force that will be needed in certain practical or engineering circumstances. If the normal force or the speed of motion gets too big, the law fails because of the excessive heat generated. It is important to realize that each of these empirical laws has its limitations, beyond which it does not really work.

That the formula $F = \mu N$ is approximately correct can be demonstrated by a simple experiment. We set up a plane, inclined at a small angle θ , and place a block of weight W on the plane. We then tilt the plane at a steeper angle, until the block just begins to slide from its own weight. The component of the weight downward along the plane is $W \sin \theta$, and this must equal the frictional force F when the block is sliding uniformly. The component of the weight normal to the plane is $W \cos \theta$, and this is the normal force N . With these values, the formula becomes $W \sin \theta = \mu W \cos \theta$, from which we get $\mu = \sin \theta / \cos \theta = \tan \theta$. If this law were exactly true, an object would start to slide at some definite inclination. If the same block is loaded by putting extra weight on it, then, although W is increased, all the forces in the formula are increased in the same proportion, and W cancels out. If μ stays constant, the loaded block will slide again at the same slope. When the angle θ is determined by trial with the original weight, it is found

that with the greater weight the block will slide at about the same angle. This will be true even when one weight is many times as great as the other, and so we conclude that the coefficient of friction is independent of the weight.

In performing this experiment it is noticeable that when the plane is tilted at about the correct angle θ , the block does not slide steadily but in a halting fashion. At one place it may stop, at another it may move with acceleration. This behavior indicates that the coefficient of friction is only roughly a constant, and varies from place to place along the plane. The same erratic behavior is observed whether the block is loaded or not. Such variations are caused by different degrees of smoothness or hardness of the plane, and perhaps dirt, oxides, or other foreign matter. The tables that list purported values of μ for "steel on steel," "copper on copper," and the like, are all false, because they ignore the factors mentioned above, which really determine μ . The friction is never due to "copper on copper," etc., but to the impurities clinging to the copper.

In experiments of the type described above, the friction is nearly independent of the velocity. Many people believe that the friction to be overcome to get something started (static friction) exceeds the force required to keep it sliding (sliding friction), but with dry metals it is very hard to show any difference. The opinion probably arises from experiences where small bits of oil or lubricant are present, or where blocks, for example, are supported by springs or other flexible supports so that they appear to bind.

It is quite difficult to do accurate quantitative experiments in friction, and the laws of friction are still not analyzed very well, in spite of the enormous engineering value of an accurate analysis. Although the law $F = \mu N$ is fairly accurate once the surfaces are standardized, the reason for this form of the law is not really understood. To show that the coefficient μ is nearly independent of velocity requires some delicate experimentation, because the apparent friction is much reduced if the lower surface vibrates very fast. When the experiment is done at very high speed, care must be taken that the objects do not vibrate relative to one another, since apparent decreases of the friction at high speed are often due to vibrations. At any rate, this friction law is another of those semiempirical laws that are not thoroughly understood, and in view of all the work that has been done it is surprising that more understanding of this phenomenon has not come about. At the present time, in fact, it is impossible even to estimate the coefficient of friction between two substances.

It was pointed out above that attempts to measure μ by sliding pure substances such as copper on copper will lead to spurious results, because the surfaces in contact are not pure copper, but are mixtures of oxides and other impurities. If we try to get absolutely pure copper, if we clean and polish the surfaces, outgas the materials in a vacuum, and take every conceivable precaution, we still do not get μ . For if we tilt the apparatus even to a vertical position, the slider will not fall off—the two pieces of copper stick together! The coefficient μ , which is ordinarily less than unity for reasonably hard surfaces, becomes several times unity! The reason for this unexpected behavior is that when the atoms in contact are all of the same kind, there is no way for the atoms to "know" that they are in different pieces of copper. When there are other atoms, in the oxides and greases and more complicated thin surface layers of contaminants in between, the atoms "know" when they are not on the same part. When we consider that it is forces between atoms that hold the copper together as a solid, it should become clear that it is impossible to get the right coefficient of friction for pure metals.

The same phenomenon can be observed in a simple home-made experiment with a flat glass plate and a glass tumbler. If the tumbler is placed on the plate and pulled along with a loop of string, it slides fairly well and one can feel the coefficient of friction; it is a little irregular, but it is a coefficient. If we now wet the glass plate and the bottom of the tumbler and pull again, we find that it binds, and if we look closely we shall find scratches, because the water is able to lift the grease and the other contaminants off the surface, and then we really have a glass-to-glass contact; this contact is so good that it holds tight and resists separation so much that the glass is torn apart; that is, it makes scratches.

12-3 Molecular forces

We shall next discuss the characteristics of molecular forces. These are forces between the atoms, and are the ultimate origin of friction. Molecular forces have never been satisfactorily explained on a basis of classical physics; it takes quantum mechanics to understand them fully. Empirically, however, the force between atoms is illustrated schematically in Fig. 12-2, where the force F between two atoms is plotted as a function of the distance r between them. There are different cases: in the water molecule, for example, the negative charges sit more on the oxygen, and the mean positions of the negative charges and of the positive charges are not at the same point; consequently, another molecule nearby feels a relatively large force, which is called a dipole-dipole force. However, for many systems the charges are very much better balanced, in particular for oxygen gas, which is perfectly symmetrical. In this case, although the minus charges and the plus charges are dispersed over the molecule, the distribution is such that the center of the minus charges and the center of the plus charges coincide. A molecule where the centers do not coincide is called a polar molecule, and charge times the separation between centers is called the dipole moment. A nonpolar molecule is one in which the centers of the charges coincide. For all nonpolar molecules, in which all the electrical forces are neutralized, it nevertheless turns out that the force at very large distances is an attraction and varies inversely as the seventh power of the distance, or $F = k/r^7$, where k is a constant that depends on the molecules. Why this is we shall learn only when we learn quantum mechanics. When there are dipoles the forces are greater. When atoms or molecules get too close they repel with a very large repulsion; that is what keeps us from falling through the floor!

These molecular forces can be demonstrated in a fairly direct way: one of these is the friction experiment with a sliding glass tumbler; another is to take two very carefully ground and lapped surfaces which are very accurately flat, so that the surfaces can be brought very close together. An example of such surfaces is the Johansson blocks that are used in machine shops as standards for making accurate length measurements. If one such block is slid over another very carefully and the upper one is lifted, the other one will adhere and also be lifted by the molecular forces, exemplifying the direct attraction between the atoms on one block for the atoms on the other block.

Nevertheless these molecular forces of attraction are still not fundamental in the sense that gravitation is fundamental; they are due to the vastly complex interactions of all the electrons and nuclei in one molecule with all the electrons and nuclei in another. Any simple-looking formula we get represents a summation of complications, so we still have not got the fundamental phenomena.

Since the molecular forces attract at large distances and repel at short distances, as shown in Fig. 12-2, we can make up solids in which all the atoms are held together by their attractions and held apart by the repulsion that sets in when they are too close together. At a certain distance d (where the graph in Fig. 12-2 crosses the axis) the forces are zero, which means that they are all balanced, so that the molecules stay that distance apart from one another. If the molecules are pushed closer together than the distance d they all show a repulsion, represented by the portion of the graph above the r -axis. To push the molecules only slightly closer together requires a great force, because the molecular repulsion rapidly becomes very great at distances less than d . If the molecules are pulled slightly apart there is a slight attraction, which increases as the separation increases. If they are pulled sufficiently hard, they will separate permanently—the bond is broken.

If the molecules are pushed only a *very small* distance closer, or pulled only a *very small* distance farther than d , the corresponding distance along the curve of Fig. 12-2 is also very small, and can then be approximated by a straight line. Therefore, in many circumstances, if the displacement is not too great *the force is proportional to the displacement*. This principle is known as Hooke's law, or the law of elasticity, which says that the force in a body which tries to restore the body

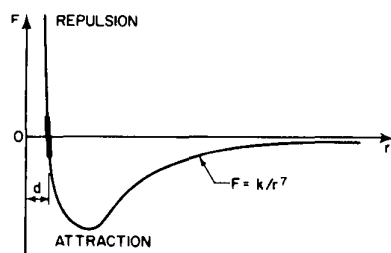


Fig. 12-2. The force between two atoms as a function of their distance of separation.

to its original condition when it is distorted is proportional to the distortion. This law, of course, holds true only if the distortion is relatively small; when it gets too large the body will be torn apart or crushed, depending on the kind of distortion. The amount of force for which Hooke's law is valid depends upon the material; for instance, for dough or putty the force is very small, but for steel it is relatively large. Hooke's law can be nicely demonstrated with a long coil spring, made of steel and suspended vertically. A suitable weight hung on the lower end of the spring produces a tiny twist throughout the length of the wire, which results in a small vertical deflection in each turn and adds up to a large displacement if there are many turns. If the total elongation produced, say, by a 100-gram weight, is measured, it is found that additional weights of 100 grams will each produce an additional elongation that is very nearly equal to the stretch that was measured for the first 100 grams. This constant ratio of force to displacement begins to change when the spring is overloaded, i.e., Hooke's law no longer holds.

12-4 Fundamental forces. Fields

We shall now discuss the only remaining forces that are fundamental. We call them fundamental in the sense that their laws are fundamentally simple. We shall first discuss electrical force. Objects carry electrical charges which consist simply of electrons or protons. If any two bodies are electrically charged, there is an electrical force between them, and if the magnitudes of the charges are q_1 and q_2 , respectively, the force varies inversely as the square of the distance between the charges, or $F = (\text{const}) q_1 q_2 / r^2$. For unlike charges, this law is like the law of gravitation, but for *like* charges the force is repulsive and the sign (direction) is reversed. The charges q_1 and q_2 can be intrinsically either positive or negative, and in any specific application of the formula the direction of the force will come out right if the q 's are given the proper plus or minus sign; the force is directed along the line between the two charges. The constant in the formula depends, of course, upon what units are used for the force, the charge, and the distance. In current practice the charge is measured in coulombs, the distance in meters, and the force in newtons. Then, in order to get the force to come out properly in newtons, the constant (which for historical reasons is written $1/4\pi\epsilon_0$) takes the numerical value

$$\epsilon_0 = 8.854 \times 10^{-12} \text{ coul}^2/\text{newton} \cdot \text{m}^2$$

or

$$1/4\pi\epsilon_0 = 8.99 \times 10^9 \text{ n} \cdot \text{m}^2/\text{coul}^2.$$

Thus the force law for static charges is

$$\mathbf{F} = q_1 q_2 \mathbf{r} / 4\pi\epsilon_0 r^3. \quad (12.2)$$

In nature, the most important charge of all is the charge on a single electron, which is 1.60×10^{-19} coulomb. In working with electrical forces between fundamental particles rather than with large charges, many people prefer the combination $(q_{e1})^2 / 4\pi\epsilon_0$, in which q_{e1} is defined as the charge on an electron. This combination occurs frequently, and to simplify calculations it has been defined by the symbol e^2 ; its numerical value in the mks system of units turns out to be $(1.52 \times 10^{-14})^2$. The advantage of using the constant in this form is that the force between two electrons in newtons can then be written simply as e^2 / r^2 , with r in meters, without all the individual constants. Electrical forces are much more complicated than this simple formula indicates, since the formula gives the force between two objects only when the objects are standing still. We shall consider the more general case shortly.

In the analysis of forces of the more fundamental kinds (not such forces as friction, but the electrical force or the gravitational force), an interesting and very important concept has been developed. Since at first sight the forces are very much more complicated than is indicated by the inverse-square laws and these laws hold true only when the interacting bodies are standing still, an improved

method is needed to deal with the very complex forces that ensue when the bodies start to move in a complicated way. Experience has shown that an approach known as the concept of a "field" is of great utility for the analysis of forces of this type. To illustrate the idea for, say, electrical force, suppose we have two electrical charges, q_1 and q_2 , located at points P and R respectively. Then the force between the charges is given by

$$\mathbf{F} = q_1 q_2 \mathbf{r} / r^3. \quad (12.3)$$

To analyze this force by means of the field concept, we say that the charge q_1 at P produces a "condition" at R , such that when the charge q_2 is placed at R it "feels" the force. This is one way, strange perhaps, of describing it; we say that the force \mathbf{F} on q_2 at R can be written in two parts. It is q_2 multiplied by a quantity \mathbf{E} that would be there whether q_2 were there or not (provided we keep all the other charges in their right places). \mathbf{E} is the "condition" produced by q_1 , we say, and \mathbf{F} is the response of q_2 to \mathbf{E} . \mathbf{E} is called an *electric field*, and it is a vector. The formula for the electric field \mathbf{E} that is produced at R by a charge q_1 at P is the charge q_1 times the constant $1/4\pi\epsilon_0$ divided by r^2 (r is the distance from P to R), and it is acting in the direction of the radius vector (the radius vector \mathbf{r} divided by its own length). The expression for \mathbf{E} is thus

$$\mathbf{E} = q_1 \mathbf{r} / 4\pi\epsilon_0 r^3. \quad (12.4)$$

We then write

$$\mathbf{F} = q_2 \mathbf{E}, \quad (12.5)$$

which expresses the force, the field, and the charge in the field. What is the point of all this? The point is to divide the analysis into two parts. One part says that something *produces* a field. The other part says that something is *acted on* by the field. By allowing us to look at the two parts independently, this separation of the analysis simplifies the calculation of a problem in many situations. If many charges are present, we first work out the total electric field produced at R by all the charges, and then, knowing the charge that is placed at R , we find the force on it.

In the case of gravitation, we can do exactly the same thing. In this case, where the force $\mathbf{F} = -Gm_1 m_2 \mathbf{r} / r^3$, we can make an analogous analysis, as follows: the force on a body in a gravitational field is the mass of that body times the field \mathbf{C} . The force on m_2 is the mass m_2 times the field \mathbf{C} produced by m_1 ; that is, $\mathbf{F} = m_2 \mathbf{C}$. Then the field \mathbf{C} produced by a body of mass m_1 is $\mathbf{C} = -Gm_1 \mathbf{r} / r^3$ and it is directed radially, as in the electrical case.

In spite of how it might at first seem, this separation of one part from another is not a triviality. It would be trivial, just another way of writing the same thing, if the laws of force were simple, but the laws of force are so complicated that it turns out that the fields have a reality that is almost independent of the objects which create them. One can do something like shake a charge and produce an effect, a field, at a distance; if one then stops moving the charge, the field keeps track of all the past, because the interaction between two particles is not instantaneous. It is desirable to have some way to remember what happened previously. If the force upon some charge depends upon where another charge was yesterday, which it does, then we need machinery to keep track of what went on yesterday, and that is the character of a field. So when the forces get more complicated, the field becomes more and more real, and this technique becomes less and less of an artificial separation.

In analyzing forces by the use of fields, we need two kinds of laws pertaining to fields. The first is the response to a field, and that gives the equations of motion. For example, the law of response of a mass to a gravitational field is that the force is equal to the mass times the gravitational field; or, if there is also a charge on the body, the response of the charge to the electric field equals the charge times the electric field. The second part of the analysis of nature in these situations is to formulate the laws which determine the strength of the field and how it is produced. These laws are sometimes called the *field equations*. We shall learn more about them in due time, but shall write down a few things about them now.

First, the most remarkable fact of all, which is true exactly and which can be easily understood, is that the total electric field produced by a number of sources is the vector sum of the electric fields produced by the first source, the second source, and so on. In other words, if we have numerous charges making a field, and if all by itself one of them would make the field E_1 , another would make the field E_2 , and so on, then we merely add the vectors to get the total field. This principle can be expressed as

$$\mathbf{E} = \mathbf{E}_1 + \mathbf{E}_2 + \mathbf{E}_3 + \cdots \quad (12.6)$$

or, in view of the definition given above,

$$\mathbf{E} = \sum_i \frac{q_i \mathbf{r}_i}{4\pi\epsilon_0 r_i^3} \quad (12.7)$$

Can the same methods be applied to gravitation? The force between two masses m_1 and m_2 was expressed by Newton as $\mathbf{F} = Gm_1m_2\mathbf{r}/r^3$. But according to the field concept, we may say that m_1 creates a field \mathbf{C} in all the surrounding space, such that the force on m_2 is given by

$$\mathbf{F} = m_2\mathbf{C} \quad (12.8)$$

By complete analogy with the electrical case,

$$\mathbf{C} = -Gm_1\mathbf{r}_1/r_1^3 \quad (12.9)$$

and the gravitational field produced by several masses is

$$\mathbf{C} = \mathbf{C}_1 + \mathbf{C}_2 + \mathbf{C}_3 + \cdots \quad (12.10)$$

In Chapter 7, in working out a case of planetary motion, we used this principle in essence. We simply added all the force vectors to get the resultant force on a planet. If we divide out the mass of the planet in question, we get Eq. (12.10).

Equations (12.6) and (12.10) express what is known as *the principle of superposition* of fields. This principle states that the total field due to all the sources is the sum of the fields due to each source. So far as we know today, for electricity this is an absolutely guaranteed law, which is true even when the force law is complicated because of the motions of the charges. There are apparent violations, but more careful analysis has always shown these to be due to the overlooking of certain moving charges. However, although the principle of superposition applies exactly for electrical forces, it is not exact for gravity if the field is too strong, and Newton's equation (12.10) is only approximate, according to Einstein's gravitational theory.

Closely related to electrical force is another kind, called magnetic force, and this too is analyzed in terms of a field. Some of the qualitative relations between electrical and magnetic forces can be illustrated by an experiment with an electron-ray tube (Fig. 12-3). At one end of such a tube is a source that emits a stream of electrons. Within the tube are arrangements for accelerating the electrons to a high speed and sending some of them in a narrow beam to a fluorescent screen at the other end of the tube. A spot of light glows in the center of the screen where the electrons strike, and this enables us to trace the electron path. On the way to the screen the electron beam passes through a narrow space between a pair of parallel metal plates, which are arranged, say, horizontally. A voltage can be applied across the plates, so that either plate can be made negative at will. When such a voltage is present, there is an electric field between the plates.

The first part of the experiment is to apply a negative voltage to the lower plate, which means that extra electrons have been placed on the lower plate. Since like charges repel, the light spot on the screen instantly shifts upward. (We could also say this in another way—that the electrons “felt” the field, and responded by deflecting upward.) We next reverse the voltage, making the *upper* plate negative. The light spot on the screen now jumps below the center, showing that the electrons in the beam were repelled by those in the plate above them. (Or we could say again

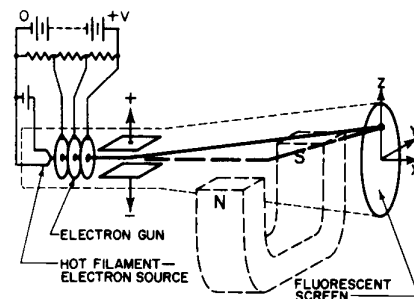


Fig. 12-3. An electron-beam tube.

that the electrons had “responded” to the field, which is now in the reverse direction.)

The second part of the experiment is to disconnect the voltage from the plates and test the effect of a magnetic field on the electron beam. This is done by means of a horseshoe magnet, whose poles are far enough apart to more or less straddle the tube. Suppose we hold the magnet below the tube in the same orientation as the letter U, with its poles up and part of the tube in between. We note that the light spot is deflected, say, upward, as the magnet approaches the tube from below. So it appears that the magnet repels the electron beam. However, it is not that simple, for if we invert the magnet without reversing the poles side-for-side, and now approach the tube from above, the spot still moves *upward*, so the electron beam is *not* repelled; instead, it appears to be attracted this time. Now we start again, restoring the magnet to its original U orientation and holding it below the tube, as before. Yes, the spot is still deflected upward; but now turn the magnet 180 degrees around a vertical axis, so that it is still in the U position but the poles are reversed side-for-side. Behold, the spot now jumps downward, and stays down, even if we invert the magnet and approach from above, as before.

To understand this peculiar behavior, we have to have a new combination of forces. We explain it thus: Across the magnet from one pole to the other there is a *magnetic field*. This field has a direction which is always away from one particular pole (which we could mark) and toward the other. Inverting the magnet did not change the direction of the field, but reversing the poles side-for-side did reverse its direction. For example, if the electron velocity were horizontal in the x -direction and the magnetic field were also horizontal but in the y -direction, the magnetic force *on the moving electrons* would be in the z -direction, i.e., up or down, depending on whether the field was in the positive or negative y -direction.

Although we shall not at the present time give the correct law of force between charges moving in an arbitrary manner, one relative to the other, because it is too complicated, we shall give one aspect of it: the complete law of the forces *if the fields are known*. The force on a charged object depends upon its motion; if, when the object is standing still at a given place, there is some force, this is taken to be proportional to the charge, the coefficient being what we call the *electric field*. When the object moves the force may be different, and the correction, the new “piece” of force, turns out to be dependent exactly *linearly on the velocity*, but at *right angles* to \mathbf{v} and to another vector quantity which we call the *magnetic induction* \mathbf{B} . If the components of the electric field \mathbf{E} and the magnetic induction \mathbf{B} are, respectively, (E_x, E_y, E_z) and (B_x, B_y, B_z) , and if the velocity \mathbf{v} has the components (v_x, v_y, v_z) , then the total electric and magnetic force on a moving charge q has the components

$$\begin{aligned} F_x &= q(E_x + v_y B_z - v_z B_y), \\ F_y &= q(E_y + v_z B_x - v_x B_z), \\ F_z &= q(E_z + v_x B_y - v_y B_x). \end{aligned} \quad (12.11)$$

If, for instance, the only component of the magnetic field were B_y and the only component of the velocity were v_x , then the only term left in the magnetic force would be a force in the z -direction, at right angles to both \mathbf{B} and \mathbf{v} .

12-5 Pseudo forces

The next kind of force we shall discuss might be called a pseudo force. In Chapter 11 we discussed the relationship between two people, Joe and Moe, who use different coordinate systems. Let us suppose that the positions of a particle as measured by Joe are x and by Moe are x' ; then the laws are as follows:

$$x = x' + s, \quad y = y', \quad z = z',$$

where s is the displacement of Moe's system relative to Joe's. If we suppose that

the laws of motion are correct for Joe, how do they look for Moe? We find first, that

$$dx/dt = dx'/dt + ds/dt.$$

Previously, we considered the case where s was constant, and we found that s made no difference in the laws of motion, since $ds/dt = 0$; ultimately, therefore, the laws of physics were the same in both systems. But another case we can take is that $s = ut$, where u is a uniform velocity in a straight line. Then s is not constant, and ds/dt is not zero, but is u , a constant. However, the acceleration d^2x/dt^2 is still the same as d^2x'/dt^2 , because $du/dt = 0$. This proves the law that we used in Chapter 10, namely, that if we move in a straight line with uniform velocity the laws of physics will look the same to us as when we are standing still. That is the Galilean transformation. But we wish to discuss the interesting case where s is still more complicated, say $s = at^2/2$. Then $ds/dt = at$ and $d^2s/dt^2 = a$, a uniform acceleration; or in a still more complicated case, the acceleration might be a function of time. This means that although the laws of force from the point of view of Joe would look like

$$m \frac{d^2x}{dt^2} = F_x,$$

the laws of force as looked upon by Moe would appear as

$$m \frac{d^2x'}{dt^2} = F_x' - ma.$$

That is, since Moe's coordinate system is accelerating with respect to Joe's, the extra term ma comes in, and Moe will have to correct his forces by that amount in order to get Newton's laws to work. In other words, here is an apparent, mysterious new force of unknown origin which arises, of course, because Moe has the wrong coordinate system. This is an example of a pseudo force; other examples occur in coordinate systems that are *rotating*.

Another example of pseudo force is what is often called "centrifugal force." An observer in a rotating coordinate system, e.g., in a rotating box, will find mysterious forces, not accounted for by any known origin of force, throwing things outward toward the walls. These forces are due merely to the fact that the observer does not have Newton's coordinate system, which is the simplest coordinate system.

Pseudo force can be illustrated by an interesting experiment in which we push a jar of water along a table, with acceleration. Gravity, of course, acts downward on the water, but because of the horizontal acceleration there is also a pseudo force acting horizontally and in a direction opposite to the acceleration. The resultant of gravity and pseudo force makes an angle with the vertical, and during the acceleration the surface of the water will be perpendicular to the resultant force, i.e., inclined at an angle with the table, with the water standing higher in the rearward side of the jar. When the push on the jar stops and the jar decelerates because of friction, the pseudo force is reversed, and the water stands higher in the forward side of the jar (Fig. 12-4).

One very important feature of pseudo forces is that they are always proportional to the masses; the same is true of gravity. The possibility exists, therefore, that *gravity itself is a pseudo force*. Is it not possible that perhaps gravitation is due simply to the fact that we do not have the right coordinate system? After all, we can always get a force proportional to the mass if we imagine that a body is accelerating. For instance, a man shut up in a box that is standing still on the earth finds himself held to the floor of the box with a certain force that is proportional to his mass. But if there were no earth at all and the box were standing still, the man inside would float in space. On the other hand, if there were no earth at all and something were *pulling* the box along with an acceleration g , then the man in the box, analyzing physics, would find a pseudo force which would pull him to the floor, just as gravity does.

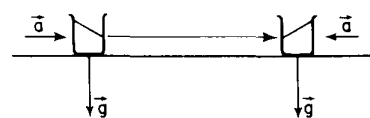


Fig. 12-4. Illustration of a pseudo force

Einstein put forward the famous hypothesis that accelerations give an imitation of gravitation, that the forces of acceleration (the pseudo forces) *cannot be distinguished* from those of gravity; it is not possible to tell how much of a given force is gravity and how much is pseudo force.

It might seem all right to consider gravity to be a pseudo force, to say that we are all held down because we are accelerating upward, but how about the people in Madagascar, on the other side of the earth—are they accelerating too? Einstein found that gravity could be considered a pseudo force only at one point at a time, and was led by his considerations to suggest that the *geometry of the world* is more complicated than ordinary Euclidean geometry. The present discussion is only qualitative, and does not pretend to convey anything more than the general idea. To give a rough idea of how gravitation could be the result of pseudo forces, we present an illustration which is purely geometrical and does not represent the real situation. Suppose that we all lived in two dimensions, and knew nothing of a third. We think we are on a plane, but suppose we are really on the surface of a sphere. And suppose that we shoot an object along the ground, with no forces on it. Where will it go? It will appear to go in a straight line, but it has to remain on the surface of a sphere, where the shortest distance between two points is along a great circle; so it goes along a great circle. If we shoot another object similarly, but in another direction, it goes along another great circle. Because we think we are on a plane, we expect that these two bodies will continue to diverge linearly with time, but careful observation will show that if they go far enough they move closer together again, as though they were attracting each other. But they are *not* attracting each other—there is just something “weird” about this geometry. This particular illustration does not describe correctly the way in which Euclid’s geometry is “weird,” but it illustrates that if we distort the geometry sufficiently it is possible that all gravitation is related in some way to pseudo forces; that is the general idea of the Einsteinian theory of gravitation.

12-6 Nuclear forces

We conclude this chapter with a brief discussion of the only other known forces, which are called *nuclear forces*. These forces are within the nuclei of atoms, and although they are much discussed, no one has ever calculated the force between two nuclei, and indeed at present there is no known law for nuclear forces. These forces have a very tiny range which is just about the same as the size of the nucleus, perhaps 10^{-13} centimeter. With particles so small and at such a tiny distance, only the quantum-mechanical laws are valid, not the Newtonian laws. In nuclear analysis we no longer think in terms of forces, and in fact we can replace the force concept with a concept of the energy of interaction of two particles, a subject that will be discussed later. Any formula that can be written for nuclear forces is a rather crude approximation which omits many complications; one might be somewhat as follows: forces within a nucleus do not vary inversely as the square of the distance, but die off exponentially over a certain distance r , as expressed by $F = (1/r^2) \exp(-r/r_0)$, where the distance r_0 is of the order of 10^{-13} centimeter. In other words, the forces disappear as soon as the particles are any great distance apart, although they are very strong within the 10^{-13} centimeter range. So far as they are understood today, the laws of nuclear force are very complex; we do not understand them in any simple way, and the whole problem of analyzing the fundamental machinery behind nuclear forces is unsolved. Attempts at a solution have led to the discovery of numerous strange particles, the π -mesons, for example, but the origin of these forces remains obscure.

Work and Potential Energy (A)

13-1 Energy of a falling body

In Chapter 4 we discussed the conservation of energy. In that discussion, we did not use Newton's laws, but it is, of course, of great interest to see how it comes about that energy is in fact conserved in accordance with these laws. For clarity we shall start with the simplest possible example, and then develop harder and harder examples.

The simplest example of the conservation of energy is a vertically falling object, one that moves only in a vertical direction. An object which changes its height under the influence of gravity alone has a kinetic energy T (or K.E.) due to its motion during the fall, and a potential energy mgh , abbreviated U or P.E., whose sum is constant:

$$\frac{1}{2}mv^2 + mgh = \text{const},$$

$\text{K E} \qquad \qquad \text{P E}$

or

$$T + U = \text{const.} \quad (13.1)$$

Now we would like to show that this statement is true. What do we mean, show it is true? From Newton's Second Law we can easily tell how the object moves, and it is easy to find out how the velocity varies with time, namely, that it increases proportionally with the time, and that the height varies as the square of the time. So if we measure the height from a zero point where the object is stationary, it is no miracle that the height turns out to be equal to the square of the velocity times a number of constants. However, let us look at it a little more closely.

Let us find out *directly* from Newton's Second Law how the kinetic energy should change, by taking the derivative of the kinetic energy with respect to time and then using Newton's laws. When we differentiate $\frac{1}{2}mv^2$ with respect to time, we obtain

$$\frac{dT}{dt} = \frac{d}{dt} (\frac{1}{2}mv^2) = \frac{1}{2}m2v \frac{dv}{dt} = mv \frac{dv}{dt}, \quad (13.2)$$

since m is assumed constant. But from Newton's Second Law, $m(dv/dt) = F$, so that

$$dT/dt = Fv. \quad (13.3)$$

In general, it will come out to be $\mathbf{F} \cdot \mathbf{v}$, but in our one-dimensional case let us leave it as the force times the velocity.

Now in our simple example the force is constant, equal to $-mg$, a vertical force (the minus sign means that it acts downward), and the velocity, of course, is the rate of change of the vertical position, or height h , with time. Thus the rate of change of the kinetic energy is $-mg(dh/dt)$, which quantity, miracle of miracles, is the rate of change of something else! It is the time rate of change of mgh ! Therefore, as time goes on, the changes in kinetic energy and in the quantity mgh are equal and opposite, so that the sum of the two quantities remains constant. Q.E.D.

We have shown, from Newton's second law of motion, that energy is conserved for constant forces when we add the potential energy mgh to the kinetic energy $\frac{1}{2}mv^2$. Now let us look into this further and see whether it can be generalized, and thus advance our understanding. Does it work only for a freely falling body, or is it more general? We expect from our discussion of the conservation of energy

13-1 Energy of a falling body

13-2 Work done by gravity

13-3 Summation of energy

13-4 Gravitational field of large objects

that it would work for an object moving from one point to another in some kind of frictionless curve, under the influence of gravity (Fig. 13-1). If the object reaches a certain height h from the original height H , then the same formula should again be right, even though the velocity is now in some direction other than the vertical. We would like to understand *why* the law is still correct. Let us follow the same analysis, finding the time rate of change of the kinetic energy. This will again be $mv(dv/dt)$, but $m(dv/dt)$ is the rate of change of the magnitude of the momentum, i.e., the *force in the direction of motion*—the tangential force F_t . Thus

$$\frac{dT}{dt} = mv \frac{dv}{dt} = F_t v.$$

Now the speed is the rate of change of distance along the curve, ds/dt , and the tangential force F_t is not mg but is weaker by the ratio of the distance ds along the path to the vertical distance dh . In other words,

$$F_t = -mg \sin \theta = -mg \frac{dh}{ds},$$

so that

$$F_t \frac{ds}{dt} = -mg \left(\frac{dh}{ds} \right) \left(\frac{ds}{dt} \right) = -mg \frac{dh}{dt},$$

since the ds 's cancel. Thus we get $-mg(dh/dt)$, which is equal to the rate of change of mgh , as before.

In order to understand exactly how the conservation of energy works in general in mechanics, we shall now discuss a number of concepts which will help us to analyze it.

First, we discuss the rate of change of kinetic energy in general in three dimensions. The kinetic energy in three dimensions is

$$T = \frac{1}{2}m(v_x^2 + v_y^2 + v_z^2).$$

When we differentiate this with respect to time, we get three terrifying terms:

$$\frac{dT}{dt} = m \left(v_x \frac{dv_x}{dt} + v_y \frac{dv_y}{dt} + v_z \frac{dv_z}{dt} \right). \quad (13.4)$$

But $m(dv_x/dt)$ is the force F_x acting on the object in the x -direction. Thus the right side of Eq. (13.4) is $F_x v_x + F_y v_y + F_z v_z$. We recall our vector analysis and recognize this as $\mathbf{F} \cdot \mathbf{v}$; therefore

$$dT/dt = \mathbf{F} \cdot \mathbf{v}. \quad (13.5)$$

This result can be derived more quickly as follows: if \mathbf{a} and \mathbf{b} are two vectors, both of which may depend upon the time, the derivative of $\mathbf{a} \cdot \mathbf{b}$ is, in general,

$$d(\mathbf{a} \cdot \mathbf{b})/dt = \mathbf{a} \cdot d\mathbf{b}/dt + (d\mathbf{a}/dt) \cdot \mathbf{b}. \quad (13.6)$$

We then use this in the form $\mathbf{a} = \mathbf{b} = \mathbf{v}$:

$$\frac{d(\frac{1}{2}mv^2)}{dt} = \frac{d(\frac{1}{2}m\mathbf{v} \cdot \mathbf{v})}{dt} = m \frac{d\mathbf{v}}{dt} \cdot \mathbf{v} = \mathbf{F} \cdot \mathbf{v} = \mathbf{F} \cdot \frac{d\mathbf{s}}{dt}. \quad (13.7)$$

Because the concepts of kinetic energy, and energy in general, are so important, various names have been given to the important terms in equations such as these. $\frac{1}{2}mv^2$ is, as we know, called *kinetic energy*. $\mathbf{F} \cdot \mathbf{v}$ is called *power*: the force acting on an object times the velocity of the object (vector "dot" product) is the power being delivered to the object by that force. We thus have a marvelous theorem: *the rate of change of kinetic energy of an object is equal to the power expended by the forces acting on it.*

However, to study the conservation of energy, we want to analyze this still more closely. Let us evaluate the change in kinetic energy in a very short time dt . If we multiply both sides of Eq. (13.7) by dt , we find that the differential change in

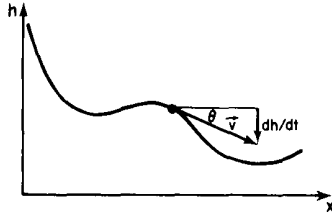


Fig. 13-1. An object moving on a frictionless curve under the influence of gravity.

the kinetic energy is the force “dot” the differential distance moved:

$$dT = \mathbf{F} \cdot d\mathbf{s}. \quad (13.8)$$

If we now integrate, we get

$$\Delta T = \int_1^2 \mathbf{F} \cdot d\mathbf{s}. \quad (13.9)$$

What does this mean? It means that if an object is moving *in any way* under the influence of a force, moving in some kind of curved path, then the change in K.E. when it goes from one point to another along the curve is equal to the integral of the component of the force along the curve times the differential displacement $d\mathbf{s}$, the integral being carried out from one point to the other. This integral also has a name; it is called the *work done by the force on the object*. We see immediately that *power equals work done per second*. We also see that it is only a component of force *in the direction of motion* that contributes to the work done. In our simple example the forces were only vertical, and had only a single component, say F_z , equal to $-mg$. No matter how the object moves in those circumstances, falling in a parabola for example, $\mathbf{F} \cdot d\mathbf{s}$, which can be written as $F_x dx + F_y dy + F_z dz$, has nothing left of it but $F_z dz = -mg dz$, because the other components of force are zero. Therefore, in our simple case,

$$\int_1^2 \mathbf{F} \cdot d\mathbf{s} = \int_{z_1}^{z_2} -mg dz = -mg(z_2 - z_1), \quad (13.10)$$

so again we find that it is only the *vertical height* from which the object falls that counts toward the potential energy.

A word about units. Since forces are measured in newtons, and we multiply by a distance in order to obtain work, work is measured in *newton · meters* (n · m), but people do not like to say newton-meters, they prefer to say *joules* (j). A newton-meter is called a joule; work is measured in joules. Power, then, is joules per second, and that is also called a *watt* (w). If we multiply watts by time, the result is the work done. The work done by the electrical company in our houses, technically, is equal to the watts times the time. That is where we get things like kilowatt hours, 1000 watts times 3600 seconds, or 3.6×10^6 joules.

Now we take another example of the law of conservation of energy. Consider an object which initially has kinetic energy and is moving very fast, and which slides against the floor with friction. It stops. At the start the kinetic energy is *not* zero, but at the end it *is* zero; there is work done by the forces, because whenever there is friction there is always a component of force in a direction opposite to that of the motion, and so energy is steadily lost. But now let us take a mass on the end of a pivot swinging in a vertical plane in a gravitational field with no friction. What happens here is different, because when the mass is going up the force is downward, and when it is coming down, the force is also downward. Thus $\mathbf{F} \cdot d\mathbf{s}$ has one sign going up and another sign coming down. At each corresponding point of the downward and upward paths the values of $\mathbf{F} \cdot d\mathbf{s}$ are exactly equal in size but of opposite sign, so the net result of the integral will be zero for this case. Thus the kinetic energy with which the mass comes back to the bottom is the same as it had when it left, that is the principle of the conservation of energy. (Note that when there are friction forces the conservation of energy seems at first sight to be invalid. We have to find another *form* of energy. It turns out, in fact, that *heat* is generated in an object when it rubs another with friction, but at the moment we supposedly do not know that.)

13-2 Work done by gravity

The next problem to be discussed is much more difficult than the above; it has to do with the case when the forces are not constant, or simply vertical, as they were in the cases we have worked out. We want to consider a planet, for example, moving around the sun, or a satellite in the space around the earth.

We shall first consider the motion of an object which starts at some point 1 and falls, say, *directly* toward the sun or toward the earth (Fig. 13-2). Will there be a law of conservation of energy in these circumstances? The only difference is that in this case, the force is *changing* as we go along, it is not just a constant. As we know, the force is GM/r^2 times the mass m , where m is the mass that moves. Now certainly when a body falls toward the earth, the kinetic energy increases as the distance fallen increases, just as it does when we do not worry about the variation of force with height. The question is whether it is possible to find another formula for potential energy different from mgh , a different function of distance away from the earth, so that conservation of energy will still be true.

This one-dimensional case is easy to treat because we know that the change in the kinetic energy is equal to the integral, from one end of the motion to the other, of $-GMm/r^2$ times the displacement dr :

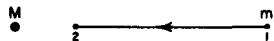


Fig. 13-2. A small mass m falls under the influence of gravity toward a large mass M .

$$T_2 - T_1 = - \int_1^2 GMm \frac{dr}{r^2}. \quad (13.11)$$

There are no cosines needed for this case because the force and the displacement are in the same direction. It is easy to integrate dr/r^2 ; the result is $-1/r$, so Eq. (13.11) becomes

$$T_2 - T_1 = +GMm \left(\frac{1}{r_2} - \frac{1}{r_1} \right). \quad (13.12)$$

Thus we have a different formula for potential energy. Equation (13.12) tells us that the quantity $(\frac{1}{2}mv^2 - GMm/r)$ calculated at point 1, at point 2, or at any other place, has a constant value.

We now have the formula for the potential energy in a gravitational field for vertical motion. Now we have an interesting problem. Can we make *perpetual motion* in a gravitational field? The gravitational field varies; in different places it is in different directions and has different strengths. Could we do something like this, using a fixed, frictionless track: start at some point and lift an object out to some other point, then move it around an arc to a third point, then lower it a certain distance, then move it in at a certain slope and pull it out some other way, so that when we bring it back to the starting point, a certain amount of work has been done by the gravitational force, and the kinetic energy of the object is increased? Can we design the curve so that it comes back moving a little bit faster than it did before, so that it goes around and around and around, and gives us perpetual motion? Since perpetual motion is impossible, we ought to find out that this is also impossible. We ought to discover the following proposition: since there is no friction the object should come back with neither higher nor lower velocity—it should be able to keep going around and around any closed path. Stated in another way, *the total work done in going around a complete cycle should be zero* for gravity forces, because if it is not zero we can get energy out by going around. (If the work turns out to be less than zero, so that we get less speed when we go around one way, then we merely go around the other way, because the forces, of course, depend only upon the position, not upon the direction; if one way is plus, the other way would be minus, so unless it is zero we will get perpetual motion by going around either way.)

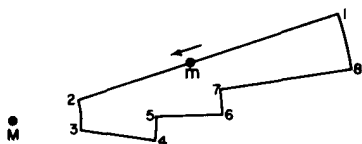


Fig. 13-3. A closed path in a gravitational field.

Is the work really zero? Let us try to demonstrate that it is. First we shall explain more or less why it is zero, and then we shall examine it a little better mathematically. Suppose that we use a simple path such as that shown in Fig. 13-3, in which a small mass is carried from point 1 to point 2, and then is made to go around a circle to 3, back to 4, then to 5, 6, 7, and 8, and finally back to 1. All of the lines are either purely radial or circular, with M as the center. How much work is done in carrying m around this path? Between points 1 and 2, it is GMm times the difference of $1/r$ between these two points:

$$W_{12} = \int_1^2 \mathbf{F} \cdot d\mathbf{s} = \int_1^2 -GMm \frac{dr}{r^2} = -GMm \left(\frac{1}{r_2} - \frac{1}{r_1} \right).$$

From 2 to 3 the force is exactly at right angles to the curve, so that $W_{23} \equiv 0$. The work from 3 to 4 is

$$W_{34} = \int_3^4 \mathbf{F} \cdot d\mathbf{s} = -GMm \left(\frac{1}{r_4} - \frac{1}{r_3} \right).$$

In the same fashion, we find that $W_{45} = 0$, $W_{56} = -GMm(1/r_6 - 1/r_5)$, $W_{67} = 0$, $W_{78} = -GMm(1/r_8 - 1/r_7)$, and $W_{81} = 0$. Thus

$$W = GMm \left(\frac{1}{r_1} - \frac{1}{r_2} + \frac{1}{r_3} - \frac{1}{r_4} + \frac{1}{r_5} - \frac{1}{r_6} + \frac{1}{r_7} - \frac{1}{r_8} \right).$$

But we note that $r_2 = r_3$, $r_4 = r_5$, $r_6 = r_7$, and $r_8 = r_1$. Therefore $W = 0$.

Of course we may wonder whether this is too trivial a curve. What if we use a *real* curve? Let us try it on a real curve. First of all, we might like to assert that a real curve could always be imitated sufficiently well by a series of sawtooth jiggles like those of Fig. 13-4, and that therefore, etc., Q.E.D., but without a little analysis, it is not obvious at first that the work done going around even a small triangle is zero. Let us magnify one of the triangles, as shown in Fig. 13-4. Is the work done in going from a to b and b to c on a triangle the same as the work done in going directly from a to c ? Suppose that the force is acting in a certain direction; let us take the triangle such that the side bc is in this direction, just as an example. We also suppose that the triangle is so small that the force is essentially constant over the entire triangle. What is the work done in going from a to c ? It is

$$W_{ac} = \int_a^c \mathbf{F} \cdot d\mathbf{s} = Fs \cos \theta,$$

since the force is constant. Now let us calculate the work done in going around the other two sides of the triangle. On the vertical side ab the force is perpendicular to ds , so that here the work is zero. On the horizontal side bc ,

$$W_{bc} = \int_b^c \mathbf{F} \cdot d\mathbf{s} = Fx.$$

Thus we see that the work done in going along the sides of a small triangle is the same as that done going on a slant, because $s \cos \theta$ is equal to x . We have proved previously that the answer is zero for any path composed of a series of notches like those of Fig. 13-3, and also that we do the same work if we cut across the corners instead of going along the notches (so long as the notches are fine enough, and we can always make them very fine); therefore, *the work done in going around any path in a gravitational field is zero*.

This is a very remarkable result. It tells us something we did not previously know about planetary motion. It tells us that when a planet moves around the sun (without any other objects around, no other forces) it moves in such a manner that the square of the speed at any point minus some constants divided by the radius at that point is always the same at every point on the orbit. For example, the closer the planet is to the sun, the faster it is going, but by how much? By the following amount: if instead of letting the planet go around the sun, we were to change the direction (but not the magnitude) of its velocity and make it move radially, and then we let it fall from some special radius to the radius of interest, the new speed would be the same as the speed it had in the actual orbit, because this is just another example of a complicated path. So long as we come back to the same distance, the kinetic energy will be the same. So, whether the motion is the real, undisturbed one, or is changed in direction by channels, by frictionless constraints, the kinetic energy with which the planet arrives at a point will be the same.

Thus, when we make a numerical analysis of the motion of the planet in its orbit, as we did earlier, we can check whether or not we are making appreciable errors by calculating this constant quantity, the energy, at every step, and it should not change. For the orbit of Table 9-2 the energy does change,* it changes by

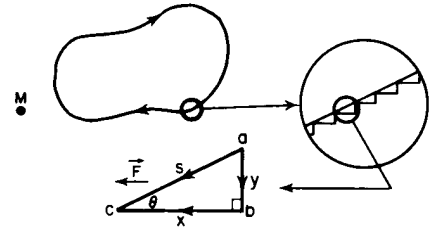


Fig. 13-4. A "smooth" closed path, showing a magnified segment of it approximated by a series of radial and circumferential steps, and an enlarged view of one step.

* The energy is $\frac{1}{2}(v_x^2 + v_y^2) - 1/r$ in the units of Table 9-2.

some 1.5 percent from the beginning to the end. Why? Either because for the numerical method we use finite intervals, or else because we made a slight mistake somewhere in arithmetic.

Let us consider the energy in another case: the problem of a mass on a spring. When we displace the mass from its balanced position, the restoring force is proportional to the displacement. In those circumstances, can we work out a law for conservation of energy? Yes, because the work done by such a force is

$$W = \int_0^x F dx = \int_0^x -kx dx = -\frac{1}{2}kx^2. \quad (13.13)$$

Therefore, for a mass on a spring we have that the kinetic energy of the oscillating mass plus $\frac{1}{2}kx^2$ is a constant. Let us see how this works. We pull the mass down; it is standing still and so its speed is zero. But x is not zero, x is at its maximum, so there is some energy, the potential energy, of course. Now we release the mass and things begin to happen (the details not to be discussed), but at any instant the kinetic plus potential energy must be a constant. For example, after the mass is on its way past the original equilibrium point, the position x equals zero, but that is when it has its biggest v^2 , and as it gets more x^2 it gets less v^2 , and so on. So the balance of x^2 and v^2 is maintained as the mass goes up and down. Thus we have another rule now, that the potential energy for a spring is $\frac{1}{2}kx^2$, if the force is $-kx$.

13-3 Summation of energy

Now we go on to the more general consideration of what happens when there are large numbers of objects. Suppose we have the complicated problem of many objects, which we label $i = 1, 2, 3, \dots$, all exerting gravitational pulls on each other. What happens then? We shall prove that if we add the kinetic energies of all the particles, and add to this the sum, over all *pairs* of particles, of their mutual gravitational potential energy, $-GMm/r_{ij}$, the total is a constant:

$$\sum_i \frac{1}{2}m_i v_i^2 + \sum_{(\text{pairs } ij)} -\frac{Gm_i m_j}{r_{ij}} = \text{const.} \quad (13.14)$$

How do we prove it? We differentiate each side with respect to time and get zero. When we differentiate $\frac{1}{2}m_i v_i^2$, we find derivatives of the velocity that are the forces, just as in Eq. (13.5). We replace these forces by the law of force that we know from Newton's law of gravity and then we notice that what is left is the same as the time derivative of

$$\sum_{\text{pairs}} -\frac{Gm_i m_j}{r_{ij}}.$$

The time derivative of the kinetic energy is

$$\begin{aligned} \frac{d}{dt} \sum_i \frac{1}{2}m_i v_i^2 &= \sum_i m_i \mathbf{v}_i \cdot \frac{d\mathbf{v}_i}{dt} \\ &= \sum_i \mathbf{F}_i \cdot \mathbf{v}_i \\ &= \sum_i \left(\sum_j -\frac{Gm_i m_j \mathbf{r}_{ij}}{r_{ij}^3} \right) \cdot \mathbf{v}_i. \end{aligned} \quad (13.15)$$

The time derivative of the potential energy is

$$\frac{d}{dt} \sum_{\text{pairs}} -\frac{Gm_i m_j}{r_{ij}} = \sum_{\text{pairs}} \left(+\frac{Gm_i m_j}{r_{ij}^2} \right) \left(\frac{dr_{ij}}{dt} \right).$$

But

$$r_{ij} = \sqrt{(x_i - x_j)^2 + (y_i - y_j)^2 + (z_i - z_j)^2},$$

so that

$$\begin{aligned}
\frac{d\mathbf{r}_{ij}}{dt} &= \frac{1}{2r_{ij}} \left[2(x_i - x_j) \left(\frac{dx_i}{dt} - \frac{dx_j}{dt} \right) \right. \\
&\quad + 2(y_i - y_j) \left(\frac{dy_i}{dt} - \frac{dy_j}{dt} \right) \\
&\quad \left. + 2(z_i - z_j) \left(\frac{dz_i}{dt} - \frac{dz_j}{dt} \right) \right] \\
&= \mathbf{r}_{ij} \cdot \frac{\mathbf{v}_i - \mathbf{v}_j}{r_{ij}} \\
&= \mathbf{r}_{ij} \cdot \frac{\mathbf{v}_i}{r_{ij}} + \mathbf{r}_{ji} \cdot \frac{\mathbf{v}_j}{r_{ji}},
\end{aligned}$$

since $\mathbf{r}_{ij} = -\mathbf{r}_{ji}$, while $r_{ij} = r_{ji}$. Thus

$$\frac{d}{dt} \sum_{\text{pairs}} -\frac{Gm_i m_j}{r_{ij}} = \sum_{\text{pairs}} \left[\frac{Gm_i m_j \mathbf{r}_{ij}}{r_{ij}^3} \cdot \mathbf{v}_i + \frac{Gm_j m_i \mathbf{r}_{ji}}{r_{ji}^3} \cdot \mathbf{v}_j \right]. \quad (13.16)$$

Now we must note carefully what $\sum_i \{\sum_j\}$ and \sum_{pairs} mean. In Eq. (13.15), $\sum_i \{\sum_j\}$ means that i takes on all values $i = 1, 2, 3, \dots$ in turn, and for each value of i , the index j takes on all values except i . Thus if $i = 3$, j takes on the values $1, 2, 4, \dots$

In Eq. (13.16), on the other hand, \sum_{pairs} means that given values of i and j occur only once. Thus the particle pair 1 and 3 contributes only one term to the sum. To keep track of this, we might agree to let i range over all values $1, 2, 3, \dots$, and for each i let j range only over values *greater* than i . Thus if $i = 3$, j could only have values $4, 5, 6, \dots$. But we notice that for each i, j value there are two contributions to the sum, one involving \mathbf{v}_i , and the other \mathbf{v}_j , and that these terms have the same appearance as those of Eq. (13.14), where *all* values of i and j (except $i = j$) are included in the sum. Therefore, by matching the terms one by one, we see that Eqs. (13.16) and (13.15) are precisely the same, but of opposite sign, so that the time derivative of the kinetic plus potential energy is indeed zero. Thus we see that, for many objects, *the kinetic energy is the sum of the contributions from each individual object*, and that the potential energy is also simple, it being also just a sum of contributions, the energies between all the pairs. We can understand *why* it should be the energy of every pair this way: Suppose that we want to find the total amount of work that must be done to bring the objects to certain distances from each other. We may do this in several steps, bringing them in from infinity where there is no force, one by one. First we bring in number one, which requires no work, since no other objects are yet present to exert force on it. Next we bring in number two, which does take some work, namely $W_{12} = -Gm_1 m_2 / r_{12}$. Now, and this is an important point, suppose we bring in the next object to position three. At any moment the force on number 3 can be written as the sum of two forces—the force exerted by number 1 and that exerted by number 2. Therefore *the work done is the sum of the works done by each*, because if \mathbf{F}_3 can be resolved into the sum of two forces,

$$\mathbf{F}_3 = \mathbf{F}_{13} + \mathbf{F}_{23},$$

then the work is

$$\int \mathbf{F}_3 \cdot d\mathbf{s} = \int \mathbf{F}_{13} \cdot d\mathbf{s} + \int \mathbf{F}_{23} \cdot d\mathbf{s} = W_{13} + W_{23}.$$

That is, the work done is the sum of the work done against the first force and the second force, as if each acted independently. Proceeding in this way, we see that the total work required to assemble the given configuration of objects is precisely the value given in Eq. (13.14) as the potential energy. It is because gravity obeys the principle of superposition of forces that we can write the potential energy as a sum over each pair of particles.

13-4 Gravitational field of large objects

Now we shall calculate the fields which are met in a few physical circumstances involving *distributions of mass*. We have not so far considered distributions of mass, only particles, so it is interesting to calculate the forces when they are produced by more than just one particle. First we shall find the gravitational force on a mass that is produced by a plane sheet of material, infinite in extent. The force on a unit mass at a given point P , produced by this sheet of material (Fig. 13-5), will of course be directed toward the sheet. Let the distance of the point from the sheet be a , and let the amount of mass per unit area of this huge sheet be μ . We shall suppose μ to be constant; it is a uniform sheet of material. Now, what small field dC is produced by the mass dm lying between ρ and $\rho + d\rho$ from the point O of the sheet nearest point P ? Answer: $dC = G(dm/r^3)$. But this field is directed along r , and we know that only the x -component of it will remain when we add all the little vector dC 's to produce C . The x -component of dC is

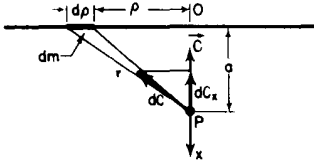


Fig. 13-5. The gravitational force F on a mass point produced by an infinite plane sheet of matter.

$$dC_x = G \frac{dmr_x}{r^3} = G \frac{dma}{r^3}.$$

Now all masses dm which are at the same distance r from P will yield the same dC_x , so we may at once write for dm the total mass in the *ring* between ρ and $\rho + d\rho$, namely $dm = \mu 2\pi\rho d\rho$ ($2\pi\rho d\rho$ is the area of a ring of radius ρ and width $d\rho$, if $d\rho \ll \rho$). Thus

$$dC_x = G\mu 2\pi\rho \frac{d\rho a}{r^3}.$$

Then, since $r^2 = \rho^2 + a^2$, $\rho d\rho = r dr$. Therefore,

$$C_x = 2\pi G\mu a \int_a^\infty \frac{dr}{r^2} = 2\pi G\mu a \left(\frac{1}{a} - \frac{1}{\infty} \right) = 2\pi G\mu. \quad (13.17)$$

Thus the force is independent of distance a ! Why? Have we made a mistake? One might think that the farther away we go, the weaker the force would be. But no! If we are close, most of the matter is pulling at an unfavorable angle; if we are far away, more of the matter is situated more favorably to exert a pull toward the plane. At any distance, the matter which is most effective lies in a certain cone. When we are farther away the force is smaller by the inverse square, but in the same cone, in the same angle, there is much *more matter*, larger by just the square of the distance! This analysis can be made rigorous by just noticing that the differential contribution in any given cone is in fact independent of the distance, because of the reciprocal variation of the strength of the force from a given mass, and the amount of mass included in the cone, with changing distance. The force is not really constant of course, because when we go on the other side of the sheet it is reversed in sign.

We have also, in effect, solved an electrical problem: if we have an electrically charged plate, with an amount σ of charge per unit area, then the electric field at a point outside the sheet is equal to $\sigma/2\epsilon_0$, and is in the outward direction if the sheet is positively charged, and inward if the sheet is negatively charged. To prove this, we merely note that G , gravity, plays the same role as $1/4\pi\epsilon_0$ for electricity.

Now suppose that we have two plates, with a positive charge $+\sigma$ on one and a negative charge $-\sigma$ on another at a distance D from the first. What is the field? Outside the two plates it is zero. Why? Because one attracts and the other repels, the force being *independent of distance*, so that the two balance out! Also, the force *between* the two plates is clearly twice as great as that from one plate, namely $E = \sigma/\epsilon_0$, and is directed from the positive plate to the negative one.

Now we come to a most interesting and important problem, whose solution we have been assuming all the time, namely, that the force produced by the earth at a point on the surface or outside it is the same as if all the mass of the earth were located at its center. The validity of this assumption is not obvious, because when we are close, some of the mass is very close to us, and some is farther away,

and so on. When we add the effects all together, it seems a miracle that the net force is exactly the same as we would get if we put all the mass in the middle!

We now demonstrate the correctness of this miracle. In order to do so, however, we shall consider a thin uniform hollow shell instead of the whole earth. Let the total mass of the shell be m , and let us calculate the *potential energy* of a particle of mass m' a distance R away from the sphere (Fig. 13-6) and show that the potential energy is the same as it would be if the mass m were a point at the center. (The potential energy is easier to work with than is the field because we do not have to worry about angles, we merely add the potential energies of all the pieces of mass.) If we call x the distance of a certain plane section from the center, then all the mass that is in a slice dx is at the same distance r from P , and the potential energy due to this ring is $-Gm' dm/r$. How much mass is in the small slice dx ? An amount

$$dm = 2\pi y \mu ds = \frac{2\pi y \mu dx}{\sin \theta} = \frac{2\pi y \mu dx a}{y} = 2\pi a \mu dx,$$

where $\mu = m/4\pi a^2$ is the surface density of mass on the spherical shell. (It is a general rule that the area of a zone of a sphere is proportional to its axial width.) Therefore the potential energy due to dm is

$$dW = -\frac{Gm' dm}{r} = -\frac{Gm' 2\pi a \mu dx}{r}.$$

But we see that

$$\begin{aligned} r^2 &= y^2 + (R - x)^2 = y^2 + x^2 + R^2 - 2Rx \\ &= a^2 + R^2 - 2Rx. \end{aligned}$$

Thus

$$2r dr = -2R dx$$

or

$$\frac{dx}{r} = \frac{dr}{R}.$$

Therefore,

$$dW = -\frac{Gm' 2\pi a \mu dr}{R},$$

and so

$$\begin{aligned} W &= -\frac{Gm' 2\pi a \mu}{R} \int_{R-a}^{R+a} dr \\ &= -\frac{Gm' 2\pi a \mu}{R} 2a = -\frac{Gm' (4\pi a^2 \mu)}{R} \\ &= -\frac{Gm' m}{R}. \end{aligned} \tag{13.18}$$

Thus, for a thin spherical shell, the potential energy of a mass m' , external to the shell, is the same as though the mass of the shell were concentrated at its center. The earth can be imagined as a series of spherical shells, each one of which contributes an energy which depends only on its mass and the distance from the center; adding them all together we get the *total mass*, and therefore the earth acts as though all the material were at the center!

But notice what happens if our point is on the *inside* of the shell. Making the same calculation, but with P on the inside, we still get the difference of the two r 's, but now in the form $a + R - (a - R) = 2R$, or twice the distance from the center. In other words, W comes out to be $W = -Gm'm/a$, which is *independent of R* and independent of position, i.e., the same energy no matter *where* we are inside. Therefore no force; no work is done when we move about inside. If the potential energy is the same no matter where an object is placed inside the sphere, there can be no force on it. So there is no force inside, there is only a force outside, and the force outside is the same as though the mass were all at the center.

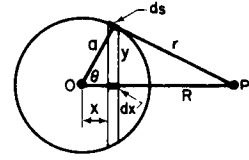


Fig. 13-6. A thin spherical shell of mass or charge.

Work and Potential Energy (conclusion)

14-1 Work

In the preceding chapter we have presented a great many new ideas and results that play a central role in physics. These ideas are so important that it seems worth while to devote a whole chapter to a closer examination of them. In the present chapter we shall not repeat the "proofs" or the specific tricks by which the results were obtained, but shall concentrate instead upon a discussion of the ideas themselves.

In learning any subject of a technical nature where mathematics plays a role, one is confronted with the task of understanding and storing away in the memory a huge body of facts and ideas, held together by certain relationships which can be "proved" or "shown" to exist between them. It is easy to confuse the proof itself with the relationship which it establishes. Clearly, the important thing to learn and to remember is the relationship, not the proof. In any particular circumstance we can either say "it can be shown that" such and such is true, or we can show it. In almost all cases, the particular proof that is used is concocted, first of all, in such form that it can be written quickly and easily on the chalkboard or on paper, and so that it will be as smooth-looking as possible. Consequently, the proof may look deceptively simple, when in fact, the author might have worked for hours trying different ways of calculating the same thing until he has found the neatest way, so as to be able to show that it can be shown in the shortest amount of time! The thing to be remembered, when seeing a proof, is not the proof itself, but rather that it *can be shown* that such and such is true. Of course, if the proof involves some mathematical procedures or "tricks" that one has not seen before, attention should be given not to the trick exactly, but to the mathematical idea involved.

It is certain that in all the demonstrations that are made in a course such as this, not one has been remembered from the time when the author studied freshman physics. Quite the contrary: he merely remembers that such and such is true, and to explain how it can be shown he invents a demonstration at the moment it is needed. Anyone who has really learned a subject should be able to follow a similar procedure, but it is no use remembering the proofs. That is why, in this chapter, we shall avoid the proofs of the various statements made previously, and merely summarize the results.

The first idea that has to be digested is *work done by a force*. The physical word "work" is not the word in the ordinary sense of "Workers of the world unite!," but is a different idea. Physical work is expressed as $\int \mathbf{F} \cdot d\mathbf{s}$, called "the line integral of \mathbf{F} dot $d\mathbf{s}$," which means that if the force, for instance, is in one direction and the object on which the force is working is displaced in a certain direction, then *only the component of force in the direction of the displacement* does any work. If, for instance, the force were constant and the displacement were a finite distance Δs , then the work done in moving the constant force through that distance is only the component of force along Δs times Δs . The rule is "force times distance," but we really mean only the component of force in the direction of the displacement times Δs or, equivalently, the component of displacement in the direction of force times F . It is evident that no work whatsoever is done by a force which is at right angles to the displacement.

Now if the vector displacement $\Delta \mathbf{s}$ is resolved into components, in other words, if the actual displacement is Δs and we want to consider it effectively as a component of displacement Δx in the x -direction, Δy in the y -direction, and Δz

14-1 Work

14-2 Constrained motion

14-3 Conservative forces

14-4 Nonconservative forces

14-5 Potentials and fields

in the z -direction, then the work done in carrying an object from one place to another can be calculated in three parts, by calculating the work done along x , along y , and along z . The work done in going along x involves only that component of force, namely F_x , and so on, so the work is $F_x \Delta x + F_y \Delta y + F_z \Delta z$. When the force is not constant, and we have a complicated curved motion, then we must resolve the path into a lot of little Δs 's, add the work done in carrying the object along each Δs , and take the limit as Δs goes to zero. This is the meaning of the "line integral."

Everything we have just said is contained in the formula $W = \int \mathbf{F} \cdot d\mathbf{s}$. It is all very well to say that it is a marvelous formula, but it is another thing to understand what it means, or what some of the consequences are.

The word "work" in physics has a meaning so different from that of the word as it is used in ordinary circumstances that it must be observed carefully that there are some peculiar circumstances in which it appears not to be the same. For example, according to the physical definition of work, if one holds a hundred-pound weight off the ground for a while, he is doing no work. Nevertheless, everyone knows that he begins to sweat, shake, and breathe harder, as if he were running up a flight of stairs. Yet running upstairs *is* considered as doing work (in running *downstairs*, one gets work out of the world, according to physics), but in simply holding an object in a fixed position, no work is done. Clearly, the physical definition of work differs from the physiological definition, for reasons we shall briefly explore.

It is a fact that when one holds a weight he has to do "physiological" work. Why should he sweat? Why should he need to consume food to hold the weight up? Why is the machinery inside him operating at full throttle, just to hold the weight up? Actually, the weight could be held up with no effort by just placing it on a table; then the table, quietly and calmly, without any supply of energy, is able to maintain the same weight at the same height! The physiological situation is something like the following. There are two kinds of muscles in the human body and in other animals: one kind, called *striated* or *skeletal* muscle, is the type of muscle we have in our arms, for example, which is under voluntary control; the other kind, called *smooth* muscle, is like the muscle in the intestines or, in the clam, the greater adductor muscle that closes the shell. The smooth muscles work very slowly, but they can hold a "set"; that is to say, if the clam tries to close its shell in a certain position, it will hold that position, even if there is a very great force trying to change it. It will hold a position under load for hours and hours without getting tired because it is very much like a table holding up a weight, it "sets" into a certain position, and the molecules just lock there temporarily with no work being done, no effort being generated by the clam. The fact that we have to generate effort to hold up a weight is simply due to the design of striated muscle. What happens is that when a nerve impulse reaches a muscle fiber, the fiber gives a little twitch and then relaxes, so that when we hold something up, enormous volleys of nerve impulses are coming in to the muscle, large numbers of twitches are maintaining the weight, while the other fibers relax. We can see this, of course: when we hold a heavy weight and get tired, we begin to shake. The reason is that the volleys are coming irregularly, and the muscle is tired and not reacting fast enough. Why such an inefficient scheme? We do not know exactly why, but evolution has not been able to develop *fast* smooth muscle. Smooth muscle would be much more effective for holding up weights because you could just stand there and it would lock in; there would be no work involved and no energy would be required. However, it has the disadvantage that it is very slow-operating.

Returning now to physics, we may ask *why* we want to calculate the work done. The answer is that it is interesting and useful to do so, since the work done on a particle by the resultant of all the forces acting on it is exactly equal to the change in kinetic energy of that particle. That is, if an object is being pushed, it picks up speed, and

$$\Delta(v^2) = \frac{2}{m} \mathbf{F} \cdot \Delta \mathbf{s}.$$

14-2 Constrained motion

Another interesting feature of forces and work is this: suppose that we have a sloping or a curved track, and a particle that must move along the track, but without friction. Or we may have a pendulum with a string and a weight; the string constrains the weight to move in a circle about the pivot point. The pivot point may be changed by having the string hit a peg, so that the path of the weight is along two circles of different radii. These are examples of what we call *fixed, frictionless constraints*.

In motion with a fixed frictionless constraint, no work is done by the constraint because the forces of constraint are always at right angles to the motion. By the "forces of constraint" we mean those forces which are applied to the object directly by the constraint itself—the contact force with the track, or the tension in the string.

The forces involved in the motion of a particle on a slope moving under the influence of gravity are quite complicated, since there is a constraint force, a gravitational force, and so on. However, if we base our calculation of the motion on conservation of energy and *the gravitational force alone*, we get the right result. This seems rather strange, because it is not strictly the right way to do it—we should use the *resultant* force. Nevertheless, the work done by the gravitational force alone will turn out to be the change in the kinetic energy, because the work done by the constraint part of the force is zero (Fig. 14-1).

The important feature here is that if a force can be analyzed as the sum of two or more "pieces" then the work done by the resultant force in going along a certain curve is the sum of the works done by the various "component" forces into which the force is analyzed. Thus if we analyze the force as being the vector sum of several effects, gravitational plus constraint forces, etc., or the x -component of all forces and the y -component of all forces, or any other way that we wish to split it up, then the work done by the net force is equal to the sum of the works done by all the parts into which we have divided the force in making the analysis.

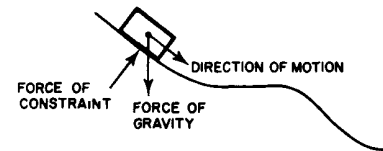


Fig. 14-1. Forces acting on a sliding body (no friction).

14-3 Conservative forces

In nature there are certain forces, that of gravity, for example, which have a very remarkable property which we call "conservative" (no political ideas involved, it is again one of those "crazy words"). If we calculate how much work is done by a force in moving an object from one point to another along some curved path, in general the work depends upon the curve, but in special cases it does not. If it does not depend upon the curve, we say that the force is a conservative force. In other words, if the integral of the force times the distance in going from position 1 to position 2 in Fig. 14-2 is calculated along curve *A* and then along *B*, we get the same number of joules, and if this is true for this pair of points on *every* curve, and if the same proposition works *no matter which pair of points* we use, then we say the force is conservative. In such circumstances, the work integral going from 1 to 2 can be evaluated in a simple manner, and we can give a formula for the result. Ordinarily it is not this easy, because we also have to specify the curve, but when we have a case where the work does not depend on the curve, then, of course, the work depends only upon the *positions* of 1 and 2.

To demonstrate this idea, consider the following. We take a "standard" point *P*, at an arbitrary location (Fig. 14-2). Then, the work line-integral from 1 to 2, which we want to calculate, can be evaluated as the work done in going from 1 to *P* plus the work done in going from *P* to 2, because the forces are conservative and the work does not depend upon the curve. Now, the work done in going from position *P* to a particular position in space is a function of that position in space. Of course it really depends on *P* also, but we hold the arbitrary point *P* fixed permanently for the analysis. If that is done, then the work done in going from point *P* to point 2 is some function of the final position of 2. It depends upon where 2 is; if we go to some other point we get a different answer.

We shall call this function of position $-U(x, y, z)$, and when we wish to refer to some particular point 2 whose coordinates are (x_2, y_2, z_2) , we shall write $U(2)$,

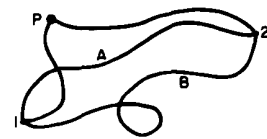


Fig. 14-2. Possible paths between two points in a field of force.

as an abbreviation for $U(x_2, y_2, z_2)$. The work done in going from point 1 to point P can be written also by going the *other way* along the integral, reversing all the ds 's. That is, the work done in going from 1 to P is *minus* the work done in going from the point P to 1:

$$\int_1^P \mathbf{F} \cdot d\mathbf{s} = \int_P^1 \mathbf{F} \cdot (-d\mathbf{s}) = - \int_P^1 \mathbf{F} \cdot d\mathbf{s}.$$

Thus the work done in going from P to 1 is $-U(1)$, and from P to 2 the work is $-U(2)$. Therefore the integral from 1 to 2 is equal to $-U(2)$ plus $[-U(1)$ backwards], or $+U(1) - U(2)$:

$$\begin{aligned} U(1) &= - \int_P^1 \mathbf{F} \cdot d\mathbf{s}, & U(2) &= - \int_P^2 \mathbf{F} \cdot d\mathbf{s}, \\ \int_1^2 \mathbf{F} \cdot d\mathbf{s} &= U(1) - U(2). \end{aligned} \quad (14.1)$$

The quantity $U(1) - U(2)$ is called the change in the potential energy, and we call U the potential energy. We shall say that when the object is located at position 2, it has potential energy $U(2)$ and at position 1 it has potential energy $U(1)$. If it is located at position P , it has zero potential energy. If we had used any other point, say Q , instead of P , it would turn out (and we shall leave it to you to demonstrate) that the *potential energy is changed only by the addition of a constant*. Since the conservation of energy depends only upon *changes*, it does not matter if we add a constant to the potential energy. Thus the point P is arbitrary.

Now, we have the following two propositions: (1) that the work done by a force is equal to the change in kinetic energy of the particle, but (2) mathematically, for a conservative force, the work done is minus the change in a function U which we call the potential energy. As a consequence of these two, we arrive at the proposition that *if only conservative forces act, the kinetic energy T plus the potential energy U remains constant*:

$$T + U = \text{constant}. \quad (14.2)$$

Let us now discuss the formulas for the potential energy for a number of cases. If we have a gravitational field that is uniform, if we are not going to heights comparable with the radius of the earth, then the force is a constant vertical force and the work done is simply the force times the vertical distance. Thus

$$U(z) = mgz, \quad (14.3)$$

and the point P which corresponds to zero potential energy happens to be any point in the plane $z = 0$. We could also have said that the potential energy is $mg(z - 6)$ if we had wanted to—all the results would, of course, be the same in our analysis except that the value of the potential energy at $z = 0$ would be $-mg6$. It makes no difference, because only *differences* in potential energy count.

The energy needed to compress a linear spring a distance x from an equilibrium point is

$$U(x) = \frac{1}{2}kx^2, \quad (14.4)$$

and the zero of potential energy is at the point $x = 0$, the equilibrium position of the spring. Again we could add any constant we wish.

The potential energy of gravitation for point masses M and m , a distance r apart, is

$$U(r) = -GMm/r. \quad (14.5)$$

The constant has been chosen here so that the potential is zero at infinity. Of course the same formula applies to electrical charges, because it is the same law:

$$U(r) = q_1q_2/4\pi\epsilon_0r. \quad (14.6)$$

Now let us actually use one of these formulas, to see whether we understand what it means. *Question*: How fast do we have to shoot a rocket away from the

earth in order for it to leave? *Solution:* The kinetic plus potential energy must be a constant; when it “leaves,” it will be millions of miles away, and if it is just barely able to leave, we may suppose that it is moving with zero speed out there, just barely going. Let a be the radius of the earth, and M its mass. The kinetic plus potential energy is then initially given by $\frac{1}{2}mv^2 - GmM/a$. At the end of the motion the two energies must be equal. The kinetic energy is taken to be zero at the end of the motion, because it is supposed to be just barely drifting away at essentially zero speed, and the potential energy is GmM divided by infinity, which is zero. So everything is zero on one side and that tells us that the square of the velocity must be $2GM/a$. But GM/a^2 is what we call the acceleration of gravity, g . Thus

$$v^2 = 2ga.$$

At what speed must a satellite travel in order to keep going around the earth? We worked this out long ago and found that $v^2 = GM/a$. Therefore to go *away* from the earth, we need $\sqrt{2}$ times the velocity we need to just go *around* the earth near its surface. We need, in other words, *twice as much energy* (because energy goes as the square of the velocity) to leave the earth as we do to go around it. Therefore the first thing that was done historically with satellites was to get one to go around the earth, which requires a speed of five miles per second. The next thing was to send a satellite away from the earth permanently; this required twice the energy, or about seven miles per second.

Now, continuing our discussion of the characteristics of potential energy, let us consider the interaction of two molecules, or two atoms, two oxygen atoms for instance. When they are very far apart, the force is one of attraction, which varies as the inverse seventh power of the distance, and when they are very close the force is a very large repulsion. If we integrate the inverse seventh power to find the work done, we find that the potential energy U , which is a function of the radial distance between the two oxygen atoms, varies as the inverse sixth power of the distance for large distances.

If we sketch the curve of the potential energy $U(r)$ as in Fig. 14-3, we thus start out at large r with an inverse sixth power, but if we come in sufficiently near we reach a point d where there is a minimum of potential energy. The minimum of potential energy at $r = d$ means this: if we start at d and move a small distance, a very small distance, the work done, which is the change in potential energy when we move this distance, is nearly zero, because there is very little change in potential energy at the bottom of the curve. Thus there is no force at this point, and so it is the equilibrium point. Another way to see that it is the equilibrium point is that it takes work to move away from d in either direction. When the two oxygen atoms have settled down, so that no more energy can be liberated from the force between them, they are in the lowest energy state, and they will be at this separation d . This is the way an oxygen molecule looks when it is cold. When we heat it up, the atoms shake and move farther apart, and we can in fact break them apart, but to do so takes a certain amount of work or energy, which is the potential energy difference between $r = d$ and $r = \infty$. When we try to push the atoms very close together the energy goes up very rapidly, because they repel each other.

The reason we bring this out is that the idea of *force* is not particularly suitable for quantum mechanics; there the idea of *energy* is most natural. We find that although forces and velocities “dissolve” and disappear when we consider the more advanced forces between nuclear matter and between molecules and so on, the energy concept remains. Therefore we find curves of potential energy in quantum mechanics books, but very rarely do we ever see a curve for the force between two molecules, because by that time people who are doing analyses are thinking in terms of energy rather than of force.

Next we note that if several conservative forces are acting on an object at the same time, then the potential energy of the object is the sum of the potential energies from each of the separate forces. This is the same proposition that we mentioned before, because if the force can be represented as a vector sum of forces, then the work done by the total force is the sum of the works done by the partial

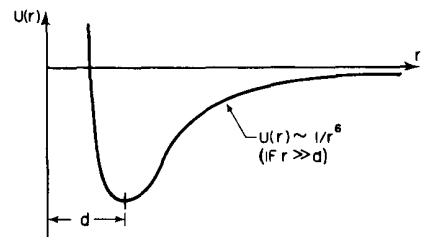


Fig. 14-3. The potential energy between two atoms as a function of the distance between them.

forces, and it can therefore be analyzed as changes in the potential energies of each of them separately. Thus the total potential energy is the sum of all the little pieces.

We could generalize this to the case of a system of many objects interacting with one another, like Jupiter, Saturn, Uranus, etc., or oxygen, nitrogen, carbon, etc., which are acting with respect to one another in pairs due to forces all of which are conservative. In these circumstances the kinetic energy in the entire system is simply the sum of the kinetic energies of all of the particular atoms or planets or whatever, and the potential energy of the system is the sum, over the pairs of particles, of the potential energy of mutual interaction of a single pair, as though the others were not there. (This is really not true for molecular forces, and the formula is somewhat more complicated; it certainly is true for Newtonian gravitation, and it is true as an approximation for molecular forces. For molecular forces there is a potential energy, but it is sometimes a more complicated function of the positions of the atoms than simply a sum of terms from pairs.) In the special case of gravity, therefore, the potential energy is the sum, over all the pairs i and j , of $-Gm_i m_j / r_{ij}$, as was indicated in Eq. (13.14). Equation (13.14) expressed mathematically the following proposition: that the total kinetic energy plus the total potential energy does not change with time. As the various planets wheel about, and turn and twist and so on, if we calculate the total kinetic energy and the total potential energy we find that the total remains constant.

14-4 Nonconservative forces

We have spent a considerable time discussing conservative forces; what about nonconservative forces? We shall take a deeper view of this than is usual, and state that there are no nonconservative forces! As a matter of fact, all the fundamental forces in nature appear to be conservative. This is not a consequence of Newton's laws. In fact, so far as Newton himself knew, the forces could be nonconservative, as friction apparently is. When we say friction *apparently* is, we are taking a modern view, in which it has been discovered that all the deep forces, the forces between the particles at the most fundamental level, are conservative.

If, for example, we analyze a system like that great globular star cluster that we saw a picture of, with the thousands of stars all interacting, then the formula for the total potential energy is simply one term plus another term, etc., summed over all pairs of stars, and the kinetic energy is the sum of the kinetic energies of all the individual stars. But the globular cluster as a whole is drifting in space too, and, if we were far enough away from it and did not see the details, could be thought of as a single object. Then if forces were applied to it, some of those forces might end up driving it forward as a whole, and we would see the center of the whole thing moving. On the other hand, some of the forces can be, so to speak, "wasted" in increasing the kinetic or potential energy of the "particles" inside. Let us suppose, for instance, that the action of these forces expands the whole cluster and makes the particles move faster. The total energy of the whole thing is really conserved, but seen from the outside with our crude eyes which cannot see the confusion of motions inside, and just thinking of the kinetic energy of the motion of the whole object as though it were a single particle, it would appear that energy is not conserved, but this is due to a lack of appreciation of what it is that we see. And that, it turns out, is the case: the total energy of the world, kinetic plus potential, is a constant when we look closely enough.

When we study matter in the finest detail at the atomic level, it is not always *easy* to separate the total energy of a thing into two parts, kinetic energy and potential energy, and such separation is not always necessary. It is *almost* always possible to do it, so let us say that it *is* always possible, and that the potential-plus-kinetic energy of the world is constant. Thus the total potential-plus-kinetic energy inside the whole world is constant, and if the "world" is a piece of isolated material, the energy is constant if there are no external forces. But as we have seen, some of the kinetic and potential energy of a thing may be internal, for instance the internal molecular motions, in the sense that we do not notice it. We know that in a glass of water everything is jiggling around, all the parts are moving

all the time, so there is a certain kinetic energy inside, which we ordinarily may not pay any attention to. We do not notice the motion of the atoms, which produces heat, and so we do not call it kinetic energy, but heat is primarily kinetic energy. Internal potential energy may also be in the form, for instance, of chemical energy: when we burn gasoline energy is liberated because the potential energies of the atoms in the new atomic arrangement are lower than in the old arrangement. It is not strictly possible to treat heat as being pure kinetic energy, for a little of the potential gets in, and vice versa for chemical energy, so we put the two together and say that the total kinetic and potential energy inside an object is partly heat, partly chemical energy, and so on. Anyway, all these different forms of internal energy are sometimes considered as "lost" energy in the sense described above; this will be made clearer when we study thermodynamics.

As another example, when friction is present it is not true that kinetic energy is lost, even though a sliding object stops and the kinetic energy seems to be lost. The kinetic energy is not lost because, of course, the atoms inside are jiggling with a greater amount of kinetic energy than before, and although we cannot see that, we can measure it by determining the temperature. Of course if we disregard the heat energy, then the conservation of energy theorem will appear to be false.

Another situation in which energy conservation appears to be false is when we study only part of a system. Naturally, the conservation of energy theorem will appear not to be true if something is interacting with something else on the outside and we neglect to take that interaction into account.

In classical physics potential energy involved only gravitation and electricity, but now we have nuclear energy and other energies also. Light, for example, would involve a new form of energy in the classical theory, but we can also, if we want to, imagine that the energy of light is the kinetic energy of a photon, and then our formula (14.2) would still be right.

14-5 Potentials and fields

We shall now discuss a few of the ideas associated with potential energy and with the idea of a *field*. Suppose we have two large objects *A* and *B* and a third very small one which is attracted gravitationally by the two, with some resultant force *F*. We have already noted in Chapter 12 that the gravitational force on a particle can be written as its mass, *m*, times another vector, *C*, which is dependent only upon the *position* of the particle:

$$\mathbf{F} = m\mathbf{C}.$$

We can analyze gravitation, then, by imagining that there is a certain vector *C* at every position in space which "acts" upon a mass which we may place there, but which is there itself whether we actually supply a mass for it to "act" on or not. *C* has three components, and each of those components is a function of (*x*, *y*, *z*), a function of position in space. Such a thing we call a *field*, and we say that the objects *A* and *B* generate the field, i.e., they "make" the vector *C*. When an object is put in a field, the force on it is equal to its mass times the value of the field vector at the point where the object is put.

We can also do the same with the potential energy. Since the potential energy, the integral of (**force**) · (*ds*) can be written as *m* times the integral of the (**field**) · (*ds*), a mere change of scale, we see that the potential energy *U*(*x*, *y*, *z*) of an object located at a point (*x*, *y*, *z*) in space can be written as *m* times another function which we may call the *potential* *Ψ*. The integral $\int \mathbf{C} \cdot d\mathbf{s} = -\Psi$, just as $\int \mathbf{F} \cdot d\mathbf{s} = -U$; there is only a scale factor between the two:

$$U = -\int \mathbf{F} \cdot d\mathbf{s} = -m \int \mathbf{C} \cdot d\mathbf{s} = m\Psi. \quad (14.7)$$

By having this function *Ψ*(*x*, *y*, *z*) at every point in space, we can immediately calculate the potential energy of an object at any point in space, namely, *U*(*x*, *y*, *z*) = *mΨ*(*x*, *y*, *z*)—rather a trivial business, it seems. But it is not really trivial, because it is sometimes much nicer to describe the field by giving the value

of Ψ everywhere in space instead of having to give \mathbf{C} . Instead of having to write three complicated components of a vector function, we can give instead the scalar function Ψ . Furthermore, it is much easier to calculate Ψ than any given component of \mathbf{C} when the field is produced by a number of masses, for since the potential is a scalar we merely add, without worrying about direction. Also, the field \mathbf{C} can be recovered easily from Ψ , as we shall shortly see. Suppose we have point masses m_1, m_2, \dots at the points 1, 2, \dots and we wish to know the potential Ψ at some arbitrary point p . This is simply the sum of the potentials at P due to the individual masses taken one by one:

$$\Psi(p) = \sum_i -\frac{Gm_i}{r_{ip}}, \quad i = 1, 2, \dots \quad (14.8)$$

In the last chapter we used this formula, that the potential is the sum of the potentials from all the different objects, to calculate the potential due to a spherical shell of matter by adding the contributions to the potential at a point from all parts of the shell. The result of this calculation is shown graphically in Fig. 14-4. It is negative, having the value zero at $r = \infty$ and varying as $1/r$ down to the radius a , and then is constant inside the shell. Outside the shell the potential is $-Gm/r$, where m is the mass of the shell, which is exactly the same as it would have been if all the mass were located at the center. But it is not *everywhere* exactly the same, for inside the shell the potential turns out to be $-Gm/a$, and is a constant! *When the potential is constant, there is no field*, or when the potential energy is constant there is no force, because if we move an object from one place to another anywhere inside the sphere the work done by the force is exactly zero. Why? Because the work done in moving the object from one place to the other is equal to minus the change in the potential energy (or, the corresponding field integral is the change of the potential). But the potential energy is the *same* at any two points inside, so there is zero change in potential energy, and therefore no work is done in going between any two points inside the shell. The only way the work can be zero for all directions of displacement is that there is no force at all.

This gives us a clue as to how we can obtain the force or the field, given the potential energy. Let us suppose that the potential energy of an object is known at the position (x, y, z) and we want to know what the force on the object is. It will not do to know the potential at only this one point, as we shall see; it requires knowledge of the potential at neighboring points as well. Why? How can we calculate the x -component of the force? (If we can do this, of course, we can also find the y - and z -components, and we will then know the whole force.) Now, if we were to move the object a small distance Δx , the work done by the force on the object would be the x -component of the force times Δx , if Δx is sufficiently small, and this should equal the change in potential energy in going from one point to the other:

$$\Delta W = -\Delta U = F_x \Delta x. \quad (14.9)$$

We have merely used the formula $\int \mathbf{F} \cdot d\mathbf{s} = -\Delta U$, but for a *very short* path. Now we divide by Δx and so find that the force is

$$F_x = -\Delta U / \Delta x. \quad (14.10)$$

Of course this is not exact. What we really want is the limit of (14.10) as Δx gets smaller and smaller, because it is only *exactly* right in the limit of infinitesimal Δx . This we recognize as the derivative of U with respect to x , and we would be inclined, therefore, to write $-dU/dx$. But U depends on x, y , and z , and the mathematicians have invented a different symbol to remind us to be very careful when we are differentiating such a function, so as to remember that we are considering that *only x varies*, and y and z do not vary. Instead of a d they simply make a "backwards 6," or ∂ . (A ∂ should have been used in the beginning of calculus because we always want to cancel that d , but we never want to cancel a ∂ !) So they write $\partial U / \partial x$, and furthermore, in moments of duress, if they want to be *very* careful, they put a line beside it with a little yz at the bottom ($\partial U / \partial x|_{yz}$),

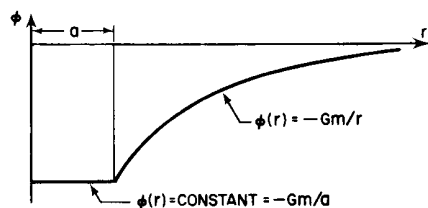


Fig. 14-4. Potential due to a spherical shell of radius a .

which means "Take the derivative of U with respect to x , keeping y and z constant." Most often we leave out the remark about what is kept constant because it is usually evident from the context, so we usually do not use the line with the y and z . However, *always* use a ∂ instead of a d as a warning that it is a derivative with some other variables kept constant. This is called a *partial derivative*; it is a derivative in which we vary only x .

Therefore, we find that the force in the x -direction is minus the partial derivative of U with respect to x :

$$F_x = -\partial U / \partial x. \quad (14.11)$$

In a similar way, the force in the y -direction can be found by differentiating U with respect to y , keeping x and z constant, and the third component, of course, is the derivative with respect to z , keeping y and x constant:

$$F_y = -\partial U / \partial y, \quad F_z = -\partial U / \partial z. \quad (14.12)$$

This is the way to get from the potential energy to the force. We get the *field* from the *potential* in exactly the same way:

$$C_x = -\partial \Psi / \partial x, \quad C_y = -\partial \Psi / \partial y, \quad C_z = -\partial \Psi / \partial z. \quad (14.13)$$

Incidentally, we shall mention here another notation, which we shall not actually use for quite a while: Since \mathbf{C} is a vector and has x -, y -, and z -components, the symbolized $\partial/\partial x$, $\partial/\partial y$, and $\partial/\partial z$ which produce the x -, y -, and z -components are something like vectors. The mathematicians have invented a glorious new symbol, ∇ , called "grad" or "gradient" which is not a quantity but an *operator* which makes a vector from a scalar. It has the following "components": The x -component of this "grad" is $\partial/\partial x$, the y -component is $\partial/\partial y$, and the z -component is $\partial/\partial z$, and then we have the fun of writing our formulas this way:

$$\mathbf{F} = -\nabla U, \quad \mathbf{C} = -\nabla \Psi. \quad (14.14)$$

Using ∇ gives us a quick way of testing whether we have a real vector equation or not, but actually Eq. (14.14) means precisely the same as Eqs. (14.11) and (14.12); it is just another way of writing them, and since we do not want to write three equations every time, we just write ∇U instead.

One more example of fields and potentials has to do with the electrical case. In the case of electricity the force on a stationary object is the charge times the electric field: $\mathbf{F} = q\mathbf{E}$. (In general, of course, the x -component of force in an electrical problem has also a part which depends on the magnetic field. It is easy to show from Eq. (12.10) that the force on a particle due to magnetic fields is always at right angles to its velocity, and also at right angles to the field. Since the force due to magnetism on a moving charge is at right angles to the velocity, *no work is done* by the magnetism on the moving charge because the motion is at right angles to the force. Therefore, in calculating theorems of kinetic energy in electric and magnetic fields we can disregard the contribution from the magnetic field, since it does not change the kinetic energy.) We suppose that there is only an electric field. Then we can calculate the energy, or work done, in the same way as for gravity, and calculate a quantity ϕ which is minus the integral of $\mathbf{E} \cdot d\mathbf{s}$, from the arbitrary fixed point to the point where we make the calculation, and then the potential energy in an electric field is just charge times this quantity ϕ :

$$\phi(\mathbf{r}) = \int \mathbf{E} \cdot d\mathbf{s},$$

$$U = q\phi.$$

Let us take, as an example, the case of two parallel metal plates, each with a surface charge of $\pm\sigma$ per unit area. This is called a parallel-plate capacitor. We found previously that there is zero force outside the plates and that there is a constant electric field between them, directed from $+$ to $-$ and of magnitude σ/ϵ_0 (Fig. 14-5). We would like to know how much work would be done in

carrying a charge from one plate to the other. The work would be the (force) \cdot (ds) integral, which can be written as charge times the potential value at plate 1 minus that at plate 2:

$$W = \int_1^2 \mathbf{F} \cdot d\mathbf{s} = q(\phi_1 - \phi_2).$$

We can actually work out the integral because the force is constant, and if we call the separation of the plates d , then the integral is easy:

$$\int_1^2 \mathbf{F} \cdot d\mathbf{s} = \frac{q\sigma}{\epsilon_0} \int_1^2 dx = \frac{q\sigma d}{\epsilon_0}.$$

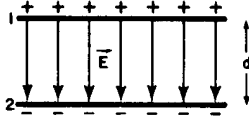


Fig. 14-5. Field between parallel plates.

The difference in potential, $\Delta\phi = \sigma d/\epsilon_0$, is called the *voltage difference*, and ϕ is measured in volts. When we say a pair of plates is charged to a certain voltage, what we mean is that the difference in electrical potential of the two plates is so-and-so many volts. For a capacitor made of two parallel plates carrying a surface charge $\pm\sigma$, the voltage, or difference in potential, of the pair of plates is $\sigma d/\epsilon_0$.

The Special Theory of Relativity

15-1 The principle of relativity

For over 200 years the equations of motion enunciated by Newton were believed to describe nature correctly, and the first time that an error in these laws was discovered, the way to correct it was also discovered. Both the error and its correction were discovered by Einstein in 1905.

Newton's Second Law, which we have expressed by the equation

$$F = d(mv)/dt,$$

was stated with the tacit assumption that m is a constant, but we now know that this is not true, and that the mass of a body increases with velocity. In Einstein's corrected formula m has the value

$$m = \frac{m_0}{\sqrt{1 - v^2/c^2}}, \quad (15.1)$$

where the "rest mass" m_0 represents the mass of a body that is not moving and c is the speed of light, which is about 3×10^5 km · sec⁻¹ or about 186,000 mi · sec⁻¹.

For those who want to learn just enough about it so they can solve problems, that is all there is to the theory of relativity—it just changes Newton's laws by introducing a correction factor to the mass. From the formula itself it is easy to see that this mass increase is very small in ordinary circumstances. If the velocity is even as great as that of a satellite, which goes around the earth at 5 mi/sec, then $v/c = 5/186,000$: putting this value into the formula shows that the correction to the mass is only one part in two to three billion, which is nearly impossible to observe. Actually, the correctness of the formula has been amply confirmed by the observation of many kinds of particles, moving at speeds ranging up to practically the speed of light. However, because the effect is ordinarily so small, it seems remarkable that it was discovered theoretically before it was discovered experimentally. Empirically, at a sufficiently high velocity, the effect is very large, but it was not discovered that way. Therefore it is interesting to see how a law that involved so delicate a modification (at the time when it was first discovered) was brought to light by a combination of experiments and physical reasoning. Contributions to the discovery were made by a number of people, the final result of whose work was Einstein's discovery.

There are really two Einstein theories of relativity. This chapter is concerned with the Special Theory of Relativity, which dates from 1905. In 1915 Einstein published an additional theory, called the General Theory of Relativity. This latter theory deals with the extension of the Special Theory to the case of the law of gravitation; we shall not discuss the General Theory here.

The principle of relativity was first stated by Newton, in one of his corollaries to the laws of motion: "The motions of bodies included in a given space are the same among themselves, whether that space is at rest or moves uniformly forward in a straight line." This means, for example, that if a space ship is drifting along at a uniform speed, all experiments performed in the space ship and all the phenomena in the space ship will appear the same as if the ship were not moving, provided, of course, that one does not look outside. That is the meaning of the principle of relativity. This is a simple enough idea, and the only question is whether it is *true* that in all experiments performed inside a moving system the laws of physics

15-1 The principle of relativity

15-2 The Lorentz transformation

15-3 The Michelson-Morley experiment

15-4 Transformation of time

15-5 The Lorentz contraction

15-6 Simultaneity

15-7 Four-vectors

15-8 Relativistic dynamics

15-9 Equivalence of mass and energy

will appear the same as they would if the system were standing still. Let us first investigate whether Newton's laws appear the same in the moving system.

Suppose that Moe is moving in the x -direction with a uniform velocity u , and he measures the position of a certain point, shown in Fig. 15-1. He designates the " x -distance" of the point in his coordinate system as x' . Joe is at rest, and measures the position of the same point, designating its x -coordinate in his system as x . The relationship of the coordinates in the two systems is clear from the diagram. After time t Moe's origin has moved a distance ut , and if the two systems originally coincided,

$$\begin{aligned}x' &= x - ut, \\y' &= y, \\z' &= z, \\t' &= t.\end{aligned}\tag{15.2}$$

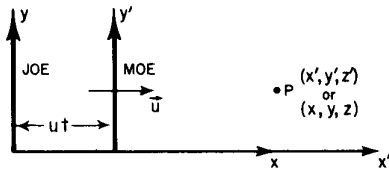


Fig. 15-1. Two coordinate systems in uniform relative motion along their x -axes.

If we substitute this transformation of coordinates into Newton's laws we find that these laws transform to the same laws in the primed system; that is, the laws of Newton are of the same form in a moving system as in a stationary system, and therefore it is impossible to tell, by making mechanical experiments, whether the system is moving or not.

The principle of relativity has been used in mechanics for a long time. It was employed by various people, in particular Huygens, to obtain the rules for the collision of billiard balls, in much the same way as we used it in Chapter 10 to discuss the conservation of momentum. In the past century interest in it was heightened as the result of investigations into the phenomena of electricity, magnetism, and light. A long series of careful studies of these phenomena by many people culminated in Maxwell's equations of the electromagnetic field, which describe electricity, magnetism, and light in one uniform system. However, the Maxwell equations did *not* seem to obey the principle of relativity. That is, if we transform Maxwell's equations by the substitution of equations 15.2, *their form does not remain the same*; therefore, in a moving space ship the electrical and optical phenomena should be different from those in a stationary ship. Thus one could use these optical phenomena to determine the speed of the ship; in particular, one could determine the absolute speed of the ship by making suitable optical or electrical measurements. One of the consequences of Maxwell's equations is that if there is a disturbance in the field such that light is generated, these electromagnetic waves go out in all directions equally and at the same speed c , or 186,000 mi/sec. Another consequence of the equations is that if the source of the disturbance is moving, the light emitted goes through space at the same speed c . This is analogous to the case of sound, the speed of sound waves being likewise independent of the motion of the source.

This independence of the motion of the source, in the case of light, brings up an interesting problem:

Suppose we are riding in a car that is going at a speed u , and light from the rear is going past the car with speed c . Differentiating the first equation in (15.2) gives

$$dx'/dt = dx/dt - u,$$

which means that according to the Galilean transformation the apparent speed of the passing light, as we measure it in the car, should not be c but should be $c - u$. For instance, if the car is going 100,000 mi/sec, and the light is going 186,000 mi/sec, then apparently the light going past the car should go 86,000 mi/sec. In any case, by measuring the speed of the light going past the car (if the Galilean transformation is correct for light), one could determine the speed of the car. A number of experiments based on this general idea were performed to determine the velocity of the earth, but they all failed—they gave *no velocity at all*. We shall discuss one of these experiments in detail, to show exactly what was done and what was the matter; something *was* the matter, of course, something was wrong with the equations of physics. What could it be?

15-2 The Lorentz transformation

When the failure of the equations of physics in the above case came to light, the first thought that occurred was that the trouble must lie in the new Maxwell equations of electrodynamics, which were only 20 years old at the time. It seemed almost obvious that these equations must be wrong, so the thing to do was to change them in such a way that under the Galilean transformation the principle of relativity would be satisfied. When this was tried, the new terms that had to be put into the equations led to predictions of new electrical phenomena that did not exist at all when tested experimentally, so this attempt had to be abandoned. Then it gradually became apparent that Maxwell's laws of electrodynamics were correct, and the trouble must be sought elsewhere.

In the meantime, H. A. Lorentz noticed a remarkable and curious thing when he made the following substitutions in the Maxwell equations:

$$\begin{aligned}x' &= \frac{x - ut}{\sqrt{1 - u^2/c^2}}, \\y' &= y, \\z' &= z, \\t' &= \frac{t - ux/c^2}{\sqrt{1 - u^2/c^2}},\end{aligned}\tag{15.3}$$

namely, Maxwell's equations remain in the same form when this transformation is applied to them! Equations (15.3) are known as a *Lorentz transformation*. Einstein, following a suggestion originally made by Poincaré, then proposed that *all the physical laws* should be of such a kind that they *remain unchanged under a Lorentz transformation*. In other words, we should change, not the laws of electrodynamics, but the laws of mechanics. How shall we change Newton's laws so that *they* will remain unchanged by the Lorentz transformation? If this goal is set, we then have to rewrite Newton's equations in such a way that the conditions we have imposed are satisfied. As it turned out, the only requirement is that the mass m in Newton's equations must be replaced by the form shown in Eq. (15.1). When this change is made, Newton's laws and the laws of electrodynamics will harmonize. Then if we use the Lorentz transformation in comparing Moe's measurements with Joe's, we shall never be able to detect whether either is moving, because the form of all the equations will be the same in both coordinate systems!

It is interesting to discuss what it means that we replace the old transformation between the coordinates and time with a new one, because the old one (Galilean) seems to be self-evident, and the new one (Lorentz) looks peculiar. We wish to know whether it is logically and experimentally possible that the new, and not the old, transformation can be correct. To find that out, it is not enough to study the laws of mechanics but, as Einstein did, we too must analyze our ideas of *space* and *time* in order to understand this transformation. We shall have to discuss these ideas and their implications for mechanics at some length, so we say in advance that the effort will be justified, since the results agree with experiment.

15-3 The Michelson-Morley experiment

As mentioned above, attempts were made to determine the absolute velocity of the earth through the hypothetical "ether" that was supposed to pervade all space. The most famous of these experiments is one performed by Michelson and Morley in 1887. It was 18 years later before the negative results of the experiment were finally explained, by Einstein.

The Michelson-Morley experiment was performed with an apparatus like that shown schematically in Fig. 15-2. This apparatus is essentially comprised of a light source A , a partially silvered glass plate B , and two mirrors C and E , all mounted on a rigid base. The mirrors are placed at equal distances L from B . The plate B splits an oncoming beam of light, and the two resulting beams con-

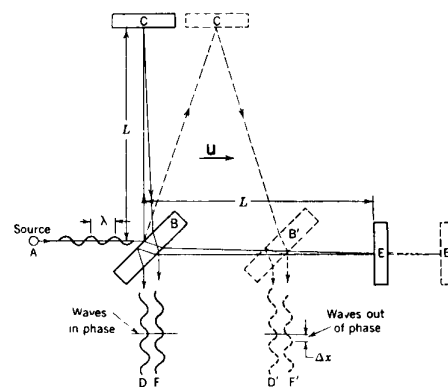


Fig. 15-2. Schematic diagram of the Michelson-Morley experiment.

tinue in mutually perpendicular directions to the mirrors, where they are reflected back to B . On arriving back at B , the two beams are recombined as two superposed beams, D and F . If the time taken for the light to go from B to E and back is the same as the time from B to C and back, the emerging beams D and F will be in phase and will reinforce each other, but if the two times differ slightly, the beams will be slightly out of phase and interference will result. If the apparatus is "at rest" in the ether, the times should be precisely equal, but if it is moving toward the right with a velocity u , there should be a difference in the times. Let us see why.

First, let us calculate the time required for the light to go from B to E and back. Let us say that the time for light to go from plate B to mirror E is t_1 , and the time for the return is t_2 . Now, while the light is on its way from B to the mirror, the apparatus moves a distance ut_1 , so the light must traverse a distance $L + ut_1$, at the speed c . We can also express this distance as ct_1 , so we have

$$ct_1 = L + ut_1, \quad \text{or} \quad t_1 = L/(c - u).$$

(This result is also obvious from the point of view that the velocity of light relative to the apparatus is $c - u$, so the time is the length L divided by $c - u$.) In a like manner, the time t_2 can be calculated. During this time the plate B advances a distance ut_2 , so the return distance of the light is $L - ut_2$. Then we have

$$ct_2 = L - ut_2, \quad \text{or} \quad t_2 = L/(c + u).$$

Then the total time is

$$t_1 + t_2 = 2Lc/(c^2 - u^2).$$

For convenience in later comparison of times we write this as

$$t_1 + t_2 = \frac{2L/c}{1 - u^2/c^2}. \quad (15.4)$$

Our second calculation will be of the time t_3 for the light to go from B to the mirror C . As before, during time t_3 the mirror C moves to the right a distance ut_3 to the position C' ; in the same time, the light travels a distance ct_3 along the hypotenuse of a triangle, which is BC' . For this right triangle we have

$$(ct_3)^2 = L^2 + (ut_3)^2$$

or

$$L^2 = c^2 t_3^2 - u^2 t_3^2 = (c^2 - u^2) t_3^2,$$

from which we get

$$t_3 = L/\sqrt{c^2 - u^2}.$$

For the return trip from C' the distance is the same, as can be seen from the symmetry of the figure; therefore the return time is also the same, and the total time is $2t_3$. With a little rearrangement of the form we can write

$$2t_3 = \frac{2L}{\sqrt{c^2 - u^2}} = \frac{2L/c}{\sqrt{1 - u^2/c^2}}. \quad (15.5)$$

We are now able to compare the times taken by the two beams of light. In expressions (15.4) and (15.5) the numerators are identical, and represent the time that would be taken if the apparatus were at rest. In the denominators, the term u^2/c^2 will be small, unless u is comparable in size to c . The denominators represent the modifications in the times caused by the motion of the apparatus. And behold, these modifications are *not the same*—the time to go to C and back is a little less than the time to E and back, even though the mirrors are equidistant from B , and all we have to do is to measure that difference with precision.

Here a minor technical point arises—suppose the two lengths L are not exactly equal? In fact, we surely cannot make them exactly equal. In that case we simply turn the apparatus 90 degrees, so that BC is in the line of motion and BE is perpendicular to the motion. Any small difference in length then becomes

unimportant, and what we look for is a *shift* in the interference fringes when we rotate the apparatus.

In carrying out the experiment, Michelson and Morley oriented the apparatus so that the line BE was nearly parallel to the earth's motion in its orbit (at certain times of the day and night). This orbital speed is about 18 miles per second, and any "ether drift" should be at least that much at some time of the day or night and at some time during the year. The apparatus was amply sensitive to observe such an effect, but no time difference was found—the velocity of the earth through the ether could not be detected. The result of the experiment was null.

The result of the Michelson-Morley experiment was very puzzling and most disturbing. The first fruitful idea for finding a way out of the impasse came from Lorentz. He suggested that material bodies contract when they are moving, and that this foreshortening is only in the direction of the motion, and also, that if the length is L_0 when a body is at rest, then when it moves with speed u parallel to its length, the new length, which we call $L_{||}$ (L -parallel), is given by

$$L_{||} = L_0 \sqrt{1 - u^2/c^2}. \quad (15.6)$$

When this modification is applied to the Michelson-Morley interferometer apparatus the distance from B to C does not change, but the distance from B to E is shortened to $L\sqrt{1 - u^2/c^2}$. Therefore Eq. (15.5) is not changed, but the L of Eq. (15.4) must be changed in accordance with Eq. (15.6). When this is done we obtain

$$t_1 + t_2 = \frac{(2L/c) \sqrt{1 - u^2/c^2}}{1 - u^2/c^2} = \frac{2L/c}{\sqrt{1 - u^2/c^2}}. \quad (15.7)$$

Comparing this result with Eq. (15.5), we see that $t_1 + t_2 = 2t_3$. So if the apparatus shrinks in the manner just described, we have a way of understanding why the Michelson-Morley experiment gives no effect at all. Although the contraction hypothesis successfully accounted for the negative result of the experiment, it was open to the objection that it was invented for the express purpose of explaining away the difficulty, and was too artificial. However, in many other experiments to discover an ether wind, similar difficulties arose, until it appeared that nature was in a "conspiracy" to thwart man by introducing some new phenomenon to undo every phenomenon that he thought would permit a measurement of u .

It was ultimately recognized, as Poincaré pointed out, that *a complete conspiracy is itself a law of nature!* Poincaré then proposed that there *is* such a law of nature, that it is not possible to discover an ether wind by *any* experiment; that is, there is no way to determine an absolute velocity.

15-4 Transformation of time

In checking out whether the contraction idea is in harmony with the facts in other experiments, it turns out that everything is correct provided that the *times* are also modified, in the manner expressed in the fourth equation of the set (15.3). That is because the time t_3 , calculated for the trip from B to C and back, is not the same when calculated by a man performing the experiment in a moving space ship as when calculated by a stationary observer who is watching the space ship. To the man in the ship the time is simply $2L/c$, but to the other observer it is $(2L/c)/\sqrt{1 - u^2/c^2}$ (Eq. 15.5). In other words, when the outsider sees the man in the space ship lighting a cigar, all the actions appear to be slower than normal, while to the man inside, everything moves at a normal rate. So not only must the lengths shorten, but also the time-measuring instruments ("clocks") must apparently slow down. That is, when the clock in the space ship records 1 second elapsed, as seen by the man in the ship, it shows $1/\sqrt{1 - u^2/c^2}$ second to the man outside.

This slowing of the clocks in a moving system is a very peculiar phenomenon, and is worth an explanation. In order to understand this, we have to watch the machinery of the clock and see what happens when it is moving. Since that is

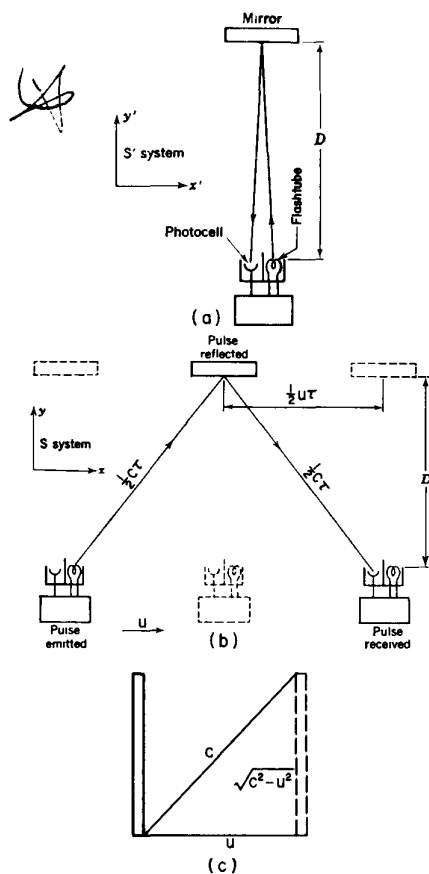


Fig. 15-3. (a) A "light clock" at rest in the S' system. (b) The same clock, moving through the S system. (c) Illustration of the diagonal path taken by the light beam in a moving "light clock."

rather difficult, we shall take a very simple kind of clock. The one we choose is rather a silly kind of clock, but it will work in principle: it is a rod (meter stick) with a mirror at each end, and when we start a light signal between the mirrors, the light keeps going up and down, making a click every time it comes down, like a standard ticking clock. We build two such clocks, with exactly the same lengths, and synchronize them by starting them together; then they agree always thereafter, because they are the same in length, and light always travels with speed c . We give one of these clocks to the man to take along in his space ship, and he mounts the rod perpendicular to the direction of motion of the ship; then the length of the rod will not change. How do we know that perpendicular lengths do not change? The men can agree to make marks on each other's y -meter stick as they pass each other. By symmetry, the two marks must come at the same y - and y' -coordinates, since otherwise, when they get together to compare results, one mark will be above or below the other, and so we could tell who was really moving.

Now let us see what happens to the moving clock. Before the man took it aboard, he agreed that it was a nice, standard clock, and when he goes along in the space ship he will not see anything peculiar. If he did, he would know he was moving—if anything at all changed because of the motion, he could tell he was moving. But the principle of relativity says this is impossible in a uniformly moving system, so nothing has changed. On the other hand, when the external observer looks at the clock going by, he sees that the light, in going from mirror to mirror, is "really" taking a zigzag path, since the rod is moving sidewise all the while. We have already analyzed such a zigzag motion in connection with the Michelson-Morley experiment. If in a given time the rod moves forward a distance proportional to u in Fig. 15-3, the distance the light travels in the same time is proportional to c , and the vertical distance is therefore proportional to $\sqrt{c^2 - u^2}$.

That is, it takes a *longer* time for light to go from end to end in the moving clock than in the stationary clock. Therefore the apparent time between clicks is longer for the moving clock, in the same proportion as shown in the hypotenuse of the triangle (that is the source of the square root expressions in our equations). From the figure it is also apparent that the greater u is, the more slowly the moving clock appears to run. Not only does this particular kind of clock run more slowly, but if the theory of relativity is correct, any other clock, operating on any principle whatsoever, would also appear to run slower, and in the same proportion—we can say this without further analysis. Why is this so?

To answer the above question, suppose we had two other clocks made exactly alike with wheels and gears, or perhaps based on radioactive decay, or something else. Then we adjust these clocks so they both run in precise synchronism with our first clocks. When light goes up and back in the first clocks and announces its arrival with a click, the new models also complete some sort of cycle, which they simultaneously announce by some doubly coincident flash, or bong, or other signal. One of these clocks is taken into the space ship, along with the first kind. Perhaps *this* clock will not run slower, but will continue to keep the same time as its stationary counterpart, and thus disagree with the other moving clock. Ah no, if that should happen, the man in the ship could use this mismatch between his two clocks to determine the speed of his ship, which we have been supposing is impossible. *We need not know anything about the machinery* of the new clock that might cause the effect—we simply know that whatever the reason, it will appear to run slow, just like the first one.

Now if *all* moving clocks run slower, if no way of measuring time gives anything but a slower rate, we shall just have to say, in a certain sense, that *time itself* appears to be slower in a space ship. All the phenomena there—the man's pulse rate, his thought processes, the time he takes to light a cigar, how long it takes to grow up and get old—all these things must be slowed down in the same proportion, because he cannot tell he is moving. The biologists and medical men sometimes say it is not quite certain that the time it takes for a cancer to develop will be longer in a space ship, but from the viewpoint of a modern physicist it is nearly certain; otherwise one could use the rate of cancer development to determine the speed of the ship!

A very interesting example of the slowing of time with motion is furnished by mu-mesons (muons), which are particles that disintegrate spontaneously after an average lifetime of 2.2×10^{-6} sec. They come to the earth in cosmic rays, and can also be produced artificially in the laboratory. Some of them disintegrate in midair, but the remainder disintegrate only after they encounter a piece of material and stop. It is clear that in its short lifetime a muon cannot travel, even at the speed of light, much more than 600 meters. But although the muons are created at the top of the atmosphere, some 10 kilometers up, yet they are actually found in a laboratory down here, in cosmic rays. How can that be? The answer is that different muons move at various speeds, some of which are very close to the speed of light. While from their own point of view they live only about $2 \mu\text{sec}$, from our point of view they live considerably longer—enough longer that they may reach the earth. The factor by which the time is increased has already been given as $1/\sqrt{1 - u^2/c^2}$. The average life has been measured quite accurately for muons of different velocities, and the values agree closely with the formula.

We do not know why the meson disintegrates or what its machinery is, but we do know its behavior satisfies the principle of relativity. That is the utility of the principle of relativity—it permits us to make predictions, even about things that otherwise we do not know much about. For example, before we have any idea at all about what makes the meson disintegrate, we can still predict that when it is moving at nine-tenths of the speed of light, the apparent length of time that it lasts is $(2.2 \times 10^{-6})/\sqrt{1 - 9^2/10^2}$ sec; and our prediction works—that is the good thing about it.

15-5 The Lorentz contraction

Now let us return to the Lorentz transformation (15.3) and try to get a better understanding of the relationship between the (x, y, z, t) and the (x', y', z', t') coordinate systems, which we shall call the S and S' systems, or Joe and Moe systems, respectively. We have already noted that the first equation is based on the Lorentz suggestion of contraction along the x -direction; how can we prove that a contraction takes place? In the Michelson-Morley experiment, we now appreciate that the *transverse* arm BC cannot change length, by the principle of relativity; yet the null result of the experiment demands that the *times* must be equal. So, in order for the experiment to give a null result, the longitudinal arm BE must appear shorter, by the square root $\sqrt{1 - u^2/c^2}$. What does this contraction mean, in terms of measurements made by Joe and Moe? Suppose that Moe, moving with the S' system in the x -direction, is measuring the x' -coordinate of some point with a meter stick. He lays the stick down x' times, so he thinks the distance is x' meters. From the viewpoint of Joe in the S system, however, Moe is using a foreshortened ruler, so the “real” distance measured is $x'\sqrt{1 - u^2/c^2}$ meters. Then if the S' system has travelled a distance ut away from the S system, the S observer would say that the same point, measured in his coordinates, is at a distance $x = x'\sqrt{1 - u^2/c^2} + ut$, or

$$x' = \frac{x - ut}{\sqrt{1 - u^2/c^2}},$$

which is the first equation of the Lorentz transformation.

15-6 Simultaneity

In an analogous way, because of the difference in time scales, the denominator expression is introduced into the fourth equation of the Lorentz transformation. The most interesting term in that equation is the ux/c^2 in the numerator, because that is quite new and unexpected. Now what does that mean? If we look at the situation carefully we see that events that occur at two separated places at the same time, as seen by Moe in S' , do *not* happen at the same time as viewed by Joe in S . If one event occurs at point x_1 at time t_0 and the other event at x_2 and t_0 (the same

time), we find that the two corresponding times t'_1 and t'_2 differ by an amount

$$t'_2 - t'_1 = \frac{u(x_1 - x_2)/c^2}{\sqrt{1 - u^2/c^2}}.$$

This circumstance is called "failure of simultaneity at a distance," and to make the idea a little clearer let us consider the following experiment.

Suppose that a man moving in a space ship (system S') has placed a clock at each end of the ship and is interested in making sure that the two clocks are in synchronism. How can the clocks be synchronized? There are many ways. One way, involving very little calculation, would be first to locate exactly the midpoint between the clocks. Then from this station we send out a light signal which will go both ways at the same speed and will arrive at both clocks, clearly, at the same time. This simultaneous arrival of the signals can be used to synchronize the clocks. Let us then suppose that the man in S' synchronizes his clocks by this particular method. Let us see whether an observer in system S would agree that the two clocks are synchronous. The man in S' has a right to believe they are, because he does not know that he is moving. But the man in S reasons that since the ship is moving forward, the clock in the front end was running away from the light signal, hence the light had to go more than halfway in order to catch up; the rear clock, however, was advancing to meet the light signal, so this distance was shorter. Therefore the signal reached the rear clock first, although the man in S' thought that the signals arrived simultaneously. We thus see that when a man in a space ship thinks the times at two locations are simultaneous, equal values of t' in his coordinate system must correspond to *different* values of t in the other coordinate system!

15-7 Four-vectors

Let us see what else we can discover in the Lorentz transformation. It is interesting to note that the transformation between the x 's and t 's is analogous in form to the transformation of the x 's and y 's that we studied in Chapter 11 for a rotation of coordinates. We then had

$$\begin{aligned} x' &= x \cos \theta + y \sin \theta, \\ y' &= y \cos \theta - x \sin \theta, \end{aligned} \quad (15.8)$$

in which the new x' mixes the old x and y , and the new y' also mixes the old x and y ; similarly, in the Lorentz transformation we find a new x' which is a mixture of x and t , and a new t' which is a mixture of t and x . So the Lorentz transformation is analogous to a rotation, only it is a "rotation" in *space and time*, which appears to be a strange concept. A check of the analogy to rotation can be made by calculating the quantity

$$x'^2 + y'^2 + z'^2 - c^2 t'^2 = x^2 + y^2 + z^2 - c^2 t^2. \quad (15.9)$$

In this equation the first three terms on each side represent, in three-dimensional geometry, the square of the distance between a point and the origin (surface of a sphere) which remains unchanged (invariant) regardless of rotation of the coordinate axes. Similarly, Eq. (15.9) shows that there is a certain combination which includes time, that is invariant to a Lorentz transformation. Thus, the analogy to a rotation is complete, and is of such a kind that vectors, i.e., quantities involving "components" which transform the same way as the coordinates and time, are also useful in connection with relativity.

Thus we contemplate an extension of the idea of vectors, which we have so far considered to have only space components, to include a time component. That is, we expect that there will be vectors with four components, three of which are like the components of an ordinary vector, and with these will be associated a fourth component, which is the analog of the time part.

This concept will be analyzed further in the next chapters, where we shall find that if the ideas of the preceding paragraph are applied to momentum, the transformation gives three space parts that are like ordinary momentum components, and a fourth component, the time part, which is the *energy*.

15-8 Relativistic dynamics

We are now ready to investigate, more generally, what form the laws of mechanics take under the Lorentz transformation. [We have thus far explained how length and time change, but not how we get the modified formula for m (Eq. 15.1). We shall do this in the next chapter.] To see the consequences of Einstein's modification of m for Newtonian mechanics, we start with the Newtonian law that force is the rate of change of momentum, or

$$\mathbf{F} = d(m\mathbf{v})/dt.$$

Momentum is still given by $m\mathbf{v}$, but when we use the new m this becomes

$$\mathbf{p} = m\mathbf{v} = \frac{m_0\mathbf{v}}{\sqrt{1 - v^2/c^2}}. \quad (15.10)$$

This is Einstein's modification of Newton's laws. Under this modification, if action and reaction are still equal (which they may not be in detail, but are in the long run), there will be conservation of momentum in the same way as before, but the quantity that is being conserved is not the old $m\mathbf{v}$ with its constant mass, but instead is the quantity shown in (15.10), which has the modified mass. When this change is made in the formula for momentum, conservation of momentum still works.

Now let us see how momentum varies with speed. In Newtonian mechanics it is proportional to the speed and, according to (15.10), over a considerable range of speed, but small compared with c , it is nearly the same in relativistic mechanics, because the square-root expression differs only slightly from 1. But when v is almost equal to c , the square-root expression approaches zero, and the momentum therefore goes toward infinity.

What happens if a constant force acts on a body for a long time? In Newtonian mechanics the body keeps picking up speed until it goes faster than light. But this is impossible in relativistic mechanics. In relativity, the body keeps picking up, not speed, but momentum, which can continually increase because the mass is increasing. After a while there is practically no acceleration in the sense of a change of velocity, but the momentum continues to increase. Of course, whenever a force produces very little change in the velocity of a body, we say that the body has a great deal of inertia, and that is exactly what our formula for relativistic mass says (see Eq. 15.10)—it says that the inertia is very great when v is nearly as great as c . As an example of this effect, to deflect the high-speed electrons in the synchrotron that is used here at Caltech, we need a magnetic field that is 2000 times stronger than would be expected on the basis of Newton's laws. In other words, the mass of the electrons in the synchrotron is 2000 times as great as their normal mass, and is as great as that of a proton! That m should be 2000 times m_0 means that $1 - v^2/c^2$ must be $1/4,000,000$, and that means that v^2/c^2 differs from 1 by one part in 4,000,000, or that v differs from c by one part in 8,000,000, so the electrons are getting pretty close to the speed of light. If the electrons and light were both to start from the synchrotron (estimated as 700 feet away) and rush out to Bridge Lab, which would arrive first? The light, of course, because light always travels faster.* How much earlier? That is too hard to tell—instead, we tell by what distance the light is ahead: it is about $1/1000$ of an inch, or $\frac{1}{4}$ the thickness of a piece of paper! When the electrons are going that fast their masses are enormous, but their speed cannot exceed the speed of light.

* The electrons would actually win the race versus *visible* light because of the index of refraction of air. A gamma ray would make out better.

Now let us look at some further consequences of relativistic change of mass. Consider the motion of the molecules in a small tank of gas. When the gas is heated, the speed of the molecules is increased, and therefore the mass is also increased and the gas is heavier. An approximate formula to express the increase of mass, for the case when the velocity is small, can be found by expanding $m_0/\sqrt{1 - v^2/c^2} = m_0(1 - v^2/c^2)^{-1/2}$ in a power series, using the binomial theorem. We get

$$m_0(1 - v^2/c^2)^{-1/2} = m_0(1 + \frac{1}{2}v^2/c^2 + \frac{3}{8}v^4/c^4 + \dots).$$

We see clearly from the formula that the series converges rapidly when v is small, and the terms after the first two or three are negligible. So we can write

$$m \cong m_0 + \frac{1}{2}m_0v^2\left(\frac{1}{c^2}\right) \quad (15.11)$$

in which the second term on the right expresses the increase of mass due to molecular velocity. When the temperature increases the v^2 increases proportionately, so we can say that the increase in mass is proportional to the increase in temperature. But since $\frac{1}{2}m_0v^2$ is the kinetic energy in the old-fashioned Newtonian sense, we can also say that the increase in mass of all this body of gas is equal to the increase in kinetic energy divided by c^2 , or $\Delta m = \Delta(\text{K.E.})/c^2$.

15-9 Equivalence of mass and energy

The above observation led Einstein to the suggestion that the mass of a body can be expressed more simply than by the formula (15.1), if we say that the mass is equal to the total energy content divided by c^2 . If Eq. (15.11) is multiplied by c^2 the result is

$$mc^2 = m_0c^2 + \frac{1}{2}m_0v^2 + \dots \quad (15.12)$$

Here, the term on the left expresses the total energy of a body, and we recognize the last term as the ordinary kinetic energy. Einstein interpreted the large constant term, m_0c^2 , to be part of the total energy of the body, an intrinsic energy known as the "rest energy."

Let us follow out the consequences of assuming, with Einstein, that *the energy of a body always equals mc^2* . As an interesting result, we shall find the formula (15.1) for the variation of mass with speed, which we have merely assumed up to now. We start with the body at rest, when its energy is m_0c^2 . Then we apply a force to the body, which starts it moving and gives it kinetic energy; therefore, since the energy has increased, the mass has increased—this is implicit in the original assumption. So long as the force continues, the energy and the mass both continue to increase. We have already seen (Chapter 13) that the rate of change of energy with time equals the force times the velocity, or

$$\frac{dE}{dt} = \mathbf{F} \cdot \mathbf{v}. \quad (15.13)$$

We also have (Chapter 9, Eq. 9.1) that $F = d(mv)/dt$. When these relations are put together with the definition of E , Eq. (15.13) becomes

$$\frac{d(mc^2)}{dt} = \mathbf{v} \cdot \frac{d(m\mathbf{v})}{dt}. \quad (15.14)$$

We wish to solve this equation for m . To do this we first use the mathematical trick of multiplying both sides by $2m$, which changes the equation to

$$c^2(2m) \frac{dm}{dt} = 2m\mathbf{v} \cdot \frac{d(m\mathbf{v})}{dt}. \quad (15.15)$$

We need to get rid of the derivatives, which can be accomplished by integrating

both sides. The quantity $(2m) dm/dt$ can be recognized as the time derivative of m^2 , and $(2mv) \cdot d(mv)/dt$ is the time derivative of $(mv)^2$. So, Eq. (15.15) is the same as

$$c^2 \frac{d(m^2)}{dt} = \frac{d(m^2 v^2)}{dt}. \quad (15.16)$$

If the derivatives of two quantities are equal, the quantities themselves differ at most by a constant, say C . This permits us to write

$$m^2 c^2 = m^2 v^2 + C. \quad (15.17)$$

We need to define the constant C more explicitly. Since Eq. (15.17) must be true for all velocities, we can choose a special case where $v = 0$, and say that in this case the mass is m_0 . Substituting these values into Eq. (15.17) gives

$$m_0^2 c^2 = 0 + C.$$

We can now use this value of C in Eq. (15.17), which becomes

$$m^2 c^2 = m^2 v^2 + m_0^2 c^2. \quad (15.18)$$

Dividing by c^2 and rearranging terms gives

$$m^2 (1 - v^2/c^2) = m_0^2,$$

from which we get

$$m = m_0 / \sqrt{1 - v^2/c^2}. \quad (15.19)$$

This is the formula (15.1), and is exactly what is necessary for the agreement between mass and energy in Eq. (15.12).

Ordinarily these energy changes represent extremely slight changes in mass, because most of the time we cannot generate much energy from a given amount of material; but in an atomic bomb of explosive energy equivalent to 20 kilotons of TNT, for example, it can be shown that the dirt after the explosion is lighter by 1 gram than the initial mass of the reacting material, because of the energy that was released, i.e., the released energy had a mass of 1 gram, according to the relationship $\Delta E = \Delta(mc^2)$. This theory of equivalence of mass and energy has been beautifully verified by experiments in which matter is annihilated—converted totally to energy: An electron and a positron come together at rest, each with a rest mass m_0 . When they come together they disintegrate and two gamma rays emerge, each with the measured energy of $m_0 c^2$. This experiment furnishes a direct determination of the energy associated with the existence of the rest mass of a particle.

Relativistic Energy and Momentum

16-1 Relativity and the philosophers

In this chapter we shall continue to discuss the principle of relativity of Einstein and Poincaré, as it affects our ideas of physics and other branches of human thought.

Poincaré made the following statement of the principle of relativity: "According to the principle of relativity, the laws of physical phenomena must be the same for a fixed observer as for an observer who has a uniform motion of translation relative to him, so that we have not, nor can we possibly have, any means of discerning whether or not we are carried along in such a motion."

When this idea descended upon the world, it caused a great stir among philosophers, particularly the "cocktail-party philosophers," who say, "Oh, it is very simple: Einstein's theory says all is relative!" In fact, a surprisingly large number of philosophers, not only those found at cocktail parties (but rather than embarrass them, we shall just call them "cocktail-party philosophers"), will say, "That all is relative is a consequence of Einstein, and it has profound influences on our ideas." In addition, they say "It has been demonstrated in physics that phenomena depend upon your frame of reference." We hear that a great deal, but it is difficult to find out what it means. Probably the frames of reference that were originally referred to were the coordinate systems which we use in the analysis of the theory of relativity. So the fact that "things depend upon your frame of reference" is supposed to have had a profound effect on modern thought. One might well wonder why, because, after all, that things depend upon one's point of view is so simple an idea that it certainly cannot have been necessary to go to all the trouble of the physical relativity theory in order to discover it. That what one sees depends upon his frame of reference is certainly known to anybody who walks around, because he sees an approaching pedestrian first from the front and then from the back; there is nothing deeper in most of the philosophy which is said to have come from the theory of relativity than the remark that "A person looks different from the front than from the back." The old story about the elephant that several blind men describe in different ways is another example, perhaps, of the theory of relativity from the philosopher's point of view.

But certainly there must be deeper things in the theory of relativity than just this simple remark that "A person looks different from the front than from the back." Of course relativity is deeper than this, because *we can make definite predictions with it*. It certainly would be rather remarkable if we could predict the behavior of nature from such a simple observation alone.

There is another school of philosophers who feel very uncomfortable about the theory of relativity, which asserts that we cannot determine our absolute velocity without looking at something outside, and who would say, "It is obvious that one cannot measure his velocity without looking outside. It is self-evident that it is *meaningless* to talk about the velocity of a thing without looking outside; the physicists are rather stupid for having thought otherwise, but it has just dawned on them that this is the case. If only we philosophers had realized what the problems were that the physicists had, we could have decided immediately by brain-work that it is impossible to tell how fast one is moving without looking outside, and we could have made an enormous contribution to physics." These philosophers are always with us, struggling in the periphery to try to tell us something, but they never really understand the subtleties and depths of the problem.

16-1 Relativity and the philosophers

16-2 The twin paradox

16-3 Transformation of velocities

16-4 Relativistic mass

16-5 Relativistic energy

Our inability to detect absolute motion is a result of *experiment* and not a result of plain thought, as we can easily illustrate. In the first place, Newton believed that it was true that one could not tell how fast he is going if he is moving with uniform velocity in a straight line. In fact, Newton first stated the principle of relativity, and one quotation made in the last chapter was a statement of Newton's. Why then did the philosophers not make all this fuss about "all is relative," or whatever, in Newton's time? Because it was not until Maxwell's theory of electrodynamics was developed that there were physical laws that suggested that one *could* measure his velocity without looking outside; soon it was found *experimentally* that one could *not*.

Now, *is* it absolutely, definitely, philosophically *necessary* that one should not be able to tell how fast he is moving without looking outside? One of the consequences of relativity was the development of a philosophy which said, "You can only define what you can measure! Since it is self-evident that one cannot measure a velocity without seeing what he is measuring it relative to, therefore it is clear that there is no *meaning* to absolute velocity. The physicists should have realized that they can talk only about what they can measure." But *that is the whole problem*: whether or not one *can define* absolute velocity is the same as the problem of whether or not one *can detect in an experiment*, without looking outside, whether he is moving. In other words, whether or not a thing is measurable is not something to be decided *a priori* by thought alone, but something that can be decided only by experiment. Given the fact that the velocity of light is 186,000 mi/sec, one will find few philosophers who will calmly state that it is self-evident that if light goes 186,000 mi/sec inside a car, and the car is going 100,000 mi/sec, that the light also goes 186,000 mi/sec past an observer on the ground. That is a shocking fact to them; the very ones who claim it is obvious find, when you give them a specific fact, that it is not obvious.

Finally, there is even a philosophy which says that one cannot detect *any* motion except by looking outside. It is simply not true in physics. True, one cannot perceive a *uniform* motion in a *straight line*, but if the whole room were *rotating* we would certainly know it, for everybody would be thrown to the wall—there would be all kinds of "centrifugal" effects. That the earth is turning on its axis can be determined without looking at the stars, by means of the so-called Foucault pendulum, for example. Therefore it is not true that "all is relative"; it is only *uniform velocity* that cannot be detected without looking outside. Uniform *rotation* about a fixed axis *can* be. When this is told to a philosopher, he is very upset that he did not really understand it, because to him it seems impossible that one should be able to determine rotation about an axis without looking outside. If the philosopher is good enough, after some time he may come back and say, "I understand. We really do not have such a thing as absolute rotation; we are really rotating *relative to the stars*, you see. And so some influence exerted by the stars on the object must cause the centrifugal force."

Now, for all we know, that is true; we have no way, at the present time, of telling whether there would have been centrifugal force if there were no stars and nebulae around. We have not been able to do the experiment of removing all the nebulae and then measuring our rotation, so we simply do not know. We must admit that the philosopher may be right. He comes back, therefore, in delight and says, "It is absolutely necessary that the world ultimately turn out to be this way: *absolute* rotation means nothing; it is only *relative* to the nebulae." Then we say to him, "Now, my friend, is it or is it not obvious that uniform velocity in a straight line, *relative to the nebulae* should produce no effects inside a car?" Now that the motion is no longer absolute, but is a motion *relative to the nebulae*, it becomes a mysterious question, and a question that can be answered only by experiment.

What, then, *are* the philosophic influences of the theory of relativity? If we limit ourselves to influences in the sense of *what kind of new ideas and suggestions* are made to the physicist by the principle of relativity, we could describe some of them as follows. The first discovery is, essentially, that even those ideas which have been held for a very long time and which have been very accurately verified might be wrong. It was a shocking discovery, of course, that Newton's laws are

wrong, after all the years in which they seemed to be accurate. Of course it is clear, not that the experiments were wrong, but that they were done over only a limited range of velocities, so small that the relativistic effects would not have been evident. But nevertheless, we now have a much more humble point of view of our physical laws—everything *can* be wrong!

Secondly, if we have a set of “strange” ideas, such as that time goes slower when one moves, and so forth, whether we *like* them or do *not* like them is an irrelevant question. The only relevant question is whether the ideas are consistent with what is found experimentally. In other words, the “strange ideas” need only agree with *experiment*, and the only reason that we have to discuss the behavior of clocks and so forth is to demonstrate that although the notion of the time dilation is strange, it is *consistent* with the way we measure time.

Finally, there is a third suggestion which is a little more technical but which has turned out to be of enormous utility in our study of other physical laws, and that is to *look at the symmetry of the laws* or, more specifically, to look for the ways in which the laws can be transformed and leave their form the same. When we discussed the theory of vectors, we noted that the fundamental laws of motion are not changed when we rotate the coordinate system, and now we learn that they are not changed when we change the space and time variables in a particular way, given by the Lorentz transformation. So this idea of studying the patterns or operations under which the fundamental laws are not changed has proved to be a very useful one.

16-2 The twin paradox

To continue our discussion of the Lorentz transformation and relativistic effects, we consider a famous so-called “paradox” of Peter and Paul, who are supposed to be twins, born at the same time. When they are old enough to drive a space ship, Paul flies away at very high speed. Because Peter, who is left on the ground, sees Paul going so fast, all of Paul’s clocks appear to go slower, his heart beats go slower, his thoughts go slower, everything goes slower, from Peter’s point of view. Of course, Paul notices nothing unusual, but if he travels around and about for a while and then comes back, he will be younger than Peter, the man on the ground! That is actually right; it is one of the consequences of the theory of relativity which has been clearly demonstrated. Just as the mu-mesons last longer when they are moving, so also will Paul last longer when he is moving. This is called a “paradox” only by the people who believe that the principle of relativity means that *all motion* is relative; they say, “Heh, heh, heh, from the point of view of Paul, can’t we say that *Peter* was moving and should therefore appear to age more slowly? By symmetry, the only possible result is that both should be the same age when they meet.” But in order for them to come back together and make the comparison, Paul must either stop at the end of the trip and make a comparison of clocks or, more simply, he has to come back, and the one who comes back must be the man who was moving, and he knows this, because he had to turn around. When he turned around, all kinds of unusual things happened in his space ship—the rockets went off, things jammed up against one wall, and so on—while Peter felt nothing.

So the way to state the rule is to say that *the man who has felt the accelerations*, who has seen things fall against the walls, and so on, is the one who would be the younger; that is the difference between them in an “absolute” sense, and it is certainly correct. When we discussed the fact that moving mu-mesons live longer, we used as an example their straight-line motion in the atmosphere. But we can also make mu-mesons in a laboratory and cause them to go in a curve with a magnet, and even under this accelerated motion, they last exactly as much longer as they do when they are moving in a straight line. Although no one has arranged an experiment explicitly so that we can get rid of the paradox, one could compare a mu-meson which is left standing with one that had gone around a complete circle, and it would surely be found that the one that went around the circle lasted longer. Although we have not actually carried out an experiment using a complete

circle, it is really not necessary, of course, because everything fits together all right. This may not satisfy those who insist that every single fact be demonstrated directly, but we confidently predict the result of the experiment in which Paul goes in a complete circle.

16-3 Transformation of velocities

The main difference between the relativity of Einstein and the relativity of Newton is that the laws of transformation connecting the coordinates and times between relatively moving systems are different. The correct transformation law, that of Lorentz, is

$$\begin{aligned}x' &= \frac{x - ut}{\sqrt{1 - u^2/c^2}}, \\y' &= y, \\z' &= z, \\t' &= \frac{t - ux/c^2}{\sqrt{1 - u^2/c^2}}.\end{aligned}\tag{16.1}$$

These equations correspond to the relatively simple case in which the relative motion of the two observers is along their common x -axes. Of course other directions of motion are possible, but the most general Lorentz transformation is rather complicated, with all four quantities mixed up together. We shall continue to use this simpler form, since it contains all the essential features of relativity.

Let us now discuss more of the consequences of this transformation. First, it is interesting to solve these equations in reverse. That is, here is a set of linear equations, four equations with four unknowns, and they can be solved in reverse, for x, y, z, t in terms of x', y', z', t' . The result is very interesting, since it tells us how a system of coordinates "at rest" looks from the point of view of one that is "moving." Of course, since the motions are relative and of uniform velocity, the man who is "moving" can say, if he wishes, that it is really the other fellow who is moving and he himself who is at rest. And since he is moving in the opposite direction, he should get the same transformation, but with the opposite sign of velocity. That is precisely what we find by manipulation, so that is consistent. If it did not come out that way, we would have real cause to worry!

$$\begin{aligned}x &= \frac{x' + ut'}{\sqrt{1 - u^2/c^2}}, \\y &= y', \\z &= z', \\t &= \frac{t' + ux'/c^2}{\sqrt{1 - u^2/c^2}}.\end{aligned}\tag{16.2}$$

Next we discuss the interesting problem of the addition of velocities in relativity. We recall that one of the original puzzles was that light travels at 186,000 mi/sec in all systems, even when they are in relative motion. This is a special case of the more general problem exemplified by the following. Suppose that an object inside a space ship is going at 100,000 mi/sec and the space ship itself is going at 100,000 mi/sec; how fast is the object inside the space ship moving from the point of view of an observer outside? We might want to say 200,000 mi/sec, which is faster than the speed of light. This is very unnerving, because it is not supposed to be going faster than the speed of light! The general problem is as follows.

Let us suppose that the object inside the ship, from the point of view of the man inside, is moving with velocity v , and that the space ship itself has a velocity u with respect to the ground. We want to know with what velocity v_x this object is moving from the point of view of the man on the ground. This is, of course, still but a special case in which the motion is in the x -direction. There will also be a

transformation for velocities in the y -direction, or for any angle; these can be worked out as needed. Inside the space ship the velocity is v_x , which means that the displacement x is equal to the velocity times the time:

$$x' = v_x t'. \quad (16.3)$$

Now we have only to calculate what the position and time are from the point of view of the outside observer for an object which has the relation (16.2) between x' and t' . So we simply substitute (16.3) into (16.2), and obtain

$$x = \frac{v_x t' + u t'}{\sqrt{1 - u^2/c^2}}. \quad (16.4)$$

But here we find x expressed in terms of t' . In order to get the velocity as seen by the man on the outside, we must divide *his distance* by *his time*, not by the *other man's time*! So we must also calculate the *time* as seen from the outside, which is

$$t = \frac{t' + u(v_x t')/c^2}{\sqrt{1 - u^2/c^2}}. \quad (16.5)$$

Now we must find the ratio of x to t , which is

$$v_x = \frac{x}{t} = \frac{u + v_x}{1 + uv_x/c^2}, \quad (16.6)$$

the square roots having cancelled. This is the law that we seek: the resultant velocity, the “summing” of two velocities, is not just the algebraic sum of two velocities (we know that it cannot be or we get in trouble), but is “corrected” by $1 + uv/c^2$.

Now let us see what happens. Suppose that you are moving inside the space ship at half the speed of light, and that the space ship itself is going at half the speed of light. Thus u is $\frac{1}{2}c$ and v is $\frac{1}{2}c$, but in the denominator uv is one-fourth, so that

$$v = \frac{\frac{1}{2}c + \frac{1}{2}c}{1 + \frac{1}{4}} = \frac{4c}{5}.$$

So, in relativity, “half” and “half” does not make “one,” it makes only “4/5.” Of course low velocities can be added quite easily in the familiar way, because so long as the velocities are small compared with the speed of light we can forget about the $(1 + uv/c^2)$ factor; but things are quite different and quite interesting at high velocity.

Let us take a limiting case. Just for fun, suppose that inside the space ship the man was observing *light itself*. In other words, $v = c$, and yet the space ship is moving. How will it look to the man on the ground? The answer will be

$$v = \frac{u + c}{1 + uc/c^2} = c \frac{u + c}{u + c} = c.$$

Therefore, if something is moving at the speed of light inside the ship, it will appear to be moving at the speed of light from the point of view of the man on the ground too! This is good, for it is, in fact, what the Einstein theory of relativity was designed to do in the first place—so it had *better* work!

Of course, there are cases in which the motion is not in the direction of the uniform translation. For example, there may be an object inside the ship which is just moving “upward” with the velocity v_y , with respect to the ship, and the ship is moving “horizontally.” Now, we simply go through the same thing, only using y 's instead of x 's, with the result

$$y = y' = v_y t',$$

so that if $v_x = 0$,

$$v_y = \frac{y}{t} = v_y \sqrt{1 - u^2/c^2}. \quad (16.7)$$

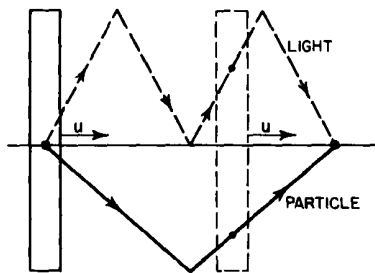


Fig 16-1. Trajectories described by a light ray and particle inside a moving clock.

Thus a sidewise velocity is no longer v_y , but $v_y\sqrt{1 - u^2/c^2}$. We found this result by substituting and combining the transformation equations, but we can also see the result directly from the principle of relativity for the following reason (it is always good to look again to see whether we can see the reason). We have already (Fig. 15-3) discussed how a possible clock might work when it is moving; the light appears to travel at an angle at the speed c in the fixed system, while it simply goes vertically with the same speed in the moving system. We found that the *vertical component* of the velocity in the fixed system is less than that of light by the factor $\sqrt{1 - u^2/c^2}$ (see Eq. 15-3). But now suppose that we let a material particle go back and forth in this same "clock," but at some integral fraction $1/n$ of the speed of light (Fig. 16-1). Then when the particle has gone back and forth once, the light will have gone exactly n times. That is, each "click" of the "particle" clock will coincide with each n th "click" of the light clock. *This fact must still be true when the whole system is moving*, because the physical phenomenon of coincidence will be a coincidence in any frame. Therefore, since the speed c_y is less than the speed of light, the speed v_y of the particle must be slower than the corresponding speed by the same square-root ratio! That is why the square root appears in any vertical velocity.

16-4 Relativistic mass

We learned in the last chapter that the mass of an object increases with velocity, but no demonstration of this was given, in the sense that we made no arguments analogous to those about the way clocks have to behave. However, we *can* show that, as a consequence of relativity plus a few other reasonable assumptions, the mass must vary in this way. (We have to say "a few other assumptions" because we cannot prove anything unless we have some laws which we assume to be true, if we expect to make meaningful deductions.) To avoid the need to study the transformation laws of force, we shall analyze a *collision*, where we need know nothing about the laws of force, except that we shall assume the conservation of momentum and energy. Also, we shall assume that the momentum of a particle which is moving is a vector and is always directed in the direction of the velocity. However, we shall not assume that the momentum is a *constant* times the velocity, as Newton did, but only that it is some *function* of velocity. We thus write the momentum vector as a certain coefficient times the vector velocity:

$$\mathbf{p} = m_v \mathbf{v}. \quad (16.8)$$

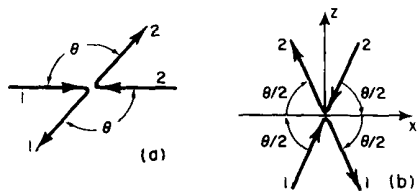


Fig. 16-2. Two views of an elastic collision between equal objects moving at the same speed in opposite directions.

We put a subscript v on the coefficient to remind us that it is a function of velocity, and we shall agree to call this coefficient m_v the "mass." Of course, when the velocity is small, it is the same mass that we would measure in the slow-moving experiments that we are used to. Now we shall try to demonstrate that the formula for m_v must be $m_0/\sqrt{1 - v^2/c^2}$, by arguing from the principle of relativity that the laws of physics must be the same in every coordinate system.

Suppose that we have two particles, like two protons, that are absolutely equal, and they are moving toward each other with exactly equal velocities. Their total momentum is zero. Now what can happen? After the collision, their directions of motion must be exactly opposite to each other, because if they are not exactly opposite, there will be a nonzero total vector momentum, and momentum would not have been conserved. Also they must have the same speeds, since they are exactly similar objects; in fact, they must have the same speed they started with, since we suppose that the energy is conserved in these collisions. So the diagram of an elastic collision, a reversible collision, will look like Fig. 16-2(a): all the arrows are the same length, all the speeds are equal. We shall suppose that such collisions can always be arranged, that any angle θ can occur, and that any speed could be used in such a collision. Next, we notice that this same collision can be viewed differently by turning the axes, and just for convenience we *shall* turn the axes, so that the horizontal splits it evenly, as in Fig. 16-2(b). It is the same collision redrawn, only with the axes turned.

Now here is the real trick: let us look at this collision from the point of view of someone riding along in a car that is moving with a speed equal to the horizontal component of the velocity of one particle. Then how does the collision look? It looks as though particle 1 is just going straight up, because it has lost its horizontal component, and it comes straight down again, also because it does not have that component. That is, the collision appears as shown in Fig. 16-3(a). Particle 2, however, was going the other way, and as we ride past it appears to fly by at some terrific speed and at a smaller angle, but we can appreciate that the angles before and after the collision are the *same*. Let us denote by u the horizontal component of the velocity of particle 2, and by w the vertical velocity of particle 1.

Now the question is, what is the vertical velocity $u \tan \alpha$? If we knew that, we could get the correct expression for the momentum, using the law of conservation of momentum in the vertical direction. Clearly, the horizontal component of the momentum is conserved: it is the same before and after the collision for both particles, and is zero for particle 1. So we need use the conservation law only for the upward velocity $u \tan \alpha$. But we *can* get the upward velocity, simply by looking at the same collision going the other way! If we look at the collision of Fig. 16-3(a) from a car to the left moving with speed u , we see the same collision, except "turned over," as shown in Fig. 16-3(b). Now particle 2 is the one that goes up and down with speed w , and particle 1 has picked up the horizontal speed u . Of course, now we *know* what the velocity $u \tan \alpha$ is: it is $w\sqrt{1 - u^2/c^2}$ (see Eq. 16.7). We know that the change in the vertical momentum of the vertically moving particle is

$$\Delta p = 2m_w w$$

(2, because it moves up and back down). The obliquely moving particle has a certain velocity v whose components we have found to be u and $w\sqrt{1 - u^2/c^2}$, and whose mass is m_v . The change in *vertical* momentum of this particle is therefore $\Delta p' = 2m_v w\sqrt{1 - u^2/c^2}$ because, in accordance with our assumed law (16.8), the momentum component is always the mass corresponding to the magnitude of the velocity times the component of the velocity in the direction of interest. Thus in order for the total momentum to be zero the vertical momenta must cancel and the ratio of the mass moving with speed v and the mass moving with speed w must therefore be

$$\frac{m_w}{m_v} = \sqrt{1 - u^2/c^2}. \quad (16.9)$$

Let us take the limiting case that w is infinitesimal. If w is very tiny indeed, it is clear that v and u are practically equal. In this case, $m_w \rightarrow m_0$ and $m_v \rightarrow m_u$. The grand result is

$$m_u = \frac{m_0}{\sqrt{1 - u^2/c^2}}. \quad (16.10)$$

It is an interesting exercise now to check whether or not Eq. (16.9) is indeed true for arbitrary values of w , assuming that Eq. (16.10) is the right formula for the mass. Note that the velocity v needed in Eq. (16.9) can be calculated from the right-angle triangle:

$$v^2 = u^2 + w^2(1 - u^2/c^2).$$

It will be found to check out automatically, although we used it only in the limit of small w .

Now, let us accept that momentum is conserved and that the mass depends upon the velocity according to (16.10) and go on to find what else we can conclude. Let us consider what is commonly called an *inelastic collision*. For simplicity, we shall suppose that two objects of the same kind, moving oppositely with equal speeds w , hit each other and stick together, to become some new, stationary object, as shown in Fig. 16-4(a). The mass m of each corresponds to w , which, as we know, is $m_0/\sqrt{1 - w^2/c^2}$. If we assume the conservation of momentum and the principle of relativity, we can demonstrate an interesting fact about the mass of the new object which has been formed. We imagine an infinitesimal velocity

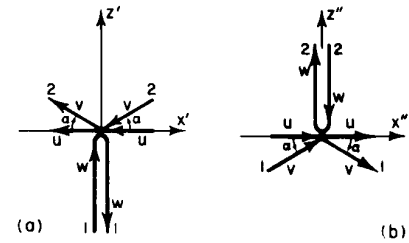


Fig. 16-3. Two more views of the collision, from moving cars.

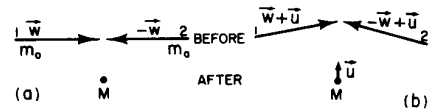


Fig. 16-4. Two views of an inelastic collision between equally massive objects.

u at right angles to w (we can do the same with finite values of u , but it is easier to understand with an infinitesimal velocity), then look at this same collision as we ride by in an elevator at the velocity $-u$. What we see is shown in Fig. 16-4(b). The composite object has an unknown mass M . Now object 1 moves with an upward component of velocity u and a horizontal component which is practically equal to w , and so also does object 2. After the collision we have the mass M moving upward with velocity u , considered very small compared with the speed of light, and also small compared with w . Momentum must be conserved, so let us estimate the momentum in the upward direction before and after the collision. Before the collision we have $p \sim 2m_w u$, and after the collision the momentum is evidently $p' = M_u u$, but M_u is essentially the same as M_0 because u is so small. These momenta must be equal because of the conservation of momentum, and therefore

$$M_0 = 2m_w. \quad (16.11)$$

The mass of the object which is formed when two equal objects collide must be twice the mass of the objects which come together. You might say, "Yes, of course, that is the conservation of mass." But not "Yes, of course," so easily, because *these masses have been enhanced* over the masses that they would be if they were standing still, yet they still contribute, to the total M , not the mass they have when standing still, but *more*. Astonishing as that may seem, in order for the conservation of momentum to work when two objects come together, the mass that they form must be greater than the rest masses of the objects, even though the objects are at rest after the collision!

16-5 Relativistic energy

In the last chapter we demonstrated that as a result of the dependence of the mass on velocity and Newton's laws, the changes in the kinetic energy of an object resulting from the total work done by the forces on it always comes out to be

$$\Delta T = (m_u - m_0)c^2 = \frac{m_0 c^2}{\sqrt{1 - u^2/c^2}} - m_0 c^2. \quad (16.12)$$

We even went further, and guessed that the total energy is the total mass times c^2 . Now we continue this discussion.

Suppose that our two equally massive objects that collide can still be "seen" inside M . For instance, a proton and a neutron are "stuck together," but are still moving about inside of M . Then, although we might at first expect the mass M to be $2m_0$, we have found that it is not $2m_0$, but $2m_w$. Since $2m_w$ is what is put in, but $2m_0$ are the rest masses of the things inside, the *excess* mass of the composite object is equal to the kinetic energy brought in. This means, of course, that *energy has inertia*. In the last chapter we discussed the heating of a gas, and showed that because the gas molecules are moving and moving things are heavier, when we put energy into the gas its molecules move faster and so the gas gets heavier. But in fact the argument is completely general, and our discussion of the inelastic collision shows that the mass is there whether or not it is *kinetic* energy. In other words, if two particles come together and produce potential or any other form of energy; if the pieces are slowed down by climbing hills, doing work against internal forces, or whatever; then it is still true that the mass is the total energy that has been put in. So we see that the conservation of mass which we have deduced above is equivalent to the conservation of energy, and therefore there is no place in the theory of relativity for strictly inelastic collisions, as there was in Newtonian mechanics. According to Newtonian mechanics it is all right for two things to collide and so form an object of mass $2m_0$ which is in no way distinct from the one that would result from putting them together slowly. Of course we know from the law of conservation of energy that there is more kinetic energy inside, but that does not affect the mass, according to Newton's laws. But now we see that this is impossible; because of the kinetic energy involved in the collision, the resulting

object will be *heavier*; therefore, it will be a *different* object. When we put the objects together gently they make something whose mass is $2m_0$; when we put them together forcefully, they make something whose mass is greater. When the mass is different, we can *tell* that it is different. So, necessarily, the conservation of energy must go along with the conservation of momentum in the theory of relativity.

This has interesting consequences. For example, suppose that we have an object whose mass M is measured, and suppose something happens so that it flies into two equal pieces moving with speed w , so that they each have a mass m_w . Now suppose that these pieces encounter enough material to slow them up until they stop; then they will have mass m_0 . How much energy will they have given to the material when they have stopped? Each will give an amount $(m_w - m_0)c^2$, by the theorem that we proved before. This much energy is left in the material in some form, as heat, potential energy, or whatever. Now $2m_w = M$, so the liberated energy is $E = (M - 2m_0)c^2$. This equation was used to estimate how much energy would be liberated under fission in the atomic bomb, for example. (Although the fragments are not exactly equal, they are nearly equal.) The mass of the uranium atom was known—it had been measured ahead of time—and the atoms into which it split, iodine, xenon, and so on, all were of known mass. By masses, we do not mean the masses while the atoms are moving, we mean the masses when the atoms are *at rest*. In other words, both M and m_0 are known. So by subtracting the two numbers one can calculate how much energy will be released if M can be made to split in “half.” For this reason poor old Einstein was called the “father” of the atomic bomb in all the newspapers. Of course, all that meant was that he could tell us ahead of time how much energy would be released if we told him what process would occur. The energy that should be liberated when an atom of uranium undergoes fission was estimated about six months before the first direct test, and as soon as the energy was in fact liberated, someone measured it directly (and if Einstein’s formula had not worked, they would have measured it anyway), and the moment they measured it they no longer needed the formula. Of course, we should not belittle Einstein, but rather should criticize the newspapers and many popular descriptions of what causes what in the history of physics and technology. The problem of how to get the thing to occur in an effective and rapid manner is a completely different matter.

The result is just as significant in chemistry. For instance, if we were to weigh the carbon dioxide molecule and compare its mass with that of the carbon and the oxygen, we could find out how much energy would be liberated when carbon and oxygen form carbon dioxide. The only trouble here is that the differences in masses are so small that it is technically very difficult to do.

Now let us turn to the question of whether we should add m_0c^2 to the kinetic energy and say from now on that the total energy of an object is mc^2 . First, if we can still *see* the component pieces of rest mass m_0 inside M , then we could say that some of the mass M of the compound object is the mechanical rest mass of the parts, part of it is kinetic energy of the parts, and part of it is potential energy of the parts. But we have discovered, in nature, particles of various kinds which undergo reactions just like the one we have treated above, in which with all the study in the world, *we cannot see the parts inside*. For instance, when a K -meson disintegrates into two pions it does so according to the law (16.11), but the idea that a K is made out of 2 π ’s is a useless idea, because it also disintegrates into 3 π ’s!

Therefore we have a *new idea*: we do not have to know what things are made of inside; we cannot and need not identify, inside a particle, which of the energy is rest energy of the parts into which it is going to disintegrate. It is not convenient and often not possible to separate the total mc^2 energy of an object into rest energy of the inside pieces, kinetic energy of the pieces, and potential energy of the pieces; instead, we simply speak of the *total energy* of the particle. We “shift the origin” of energy by adding a constant m_0c^2 to everything, and say that the total energy of a particle is the mass in motion times c^2 , and when the object is standing still, the energy is the mass at rest times c^2 .

Finally, we find that the velocity v , momentum P , and total energy E are related in a rather simple way. That the mass in motion at speed v is the mass m_0 at rest divided by $\sqrt{1 - v^2/c^2}$, surprisingly enough, is rarely used. Instead, the following relations are easily proved, and turn out to be very useful:

$$E^2 - P^2c^2 = m_0^2c^4 \quad (16.13)$$

and

$$Pc = Ev/c. \quad (16.14)$$

Space-Time

17-1 The geometry of space-time

The theory of relativity shows us that the relationships of positions and times as measured in one coordinate system and another are not what we would have expected on the basis of our intuitive ideas. It is very important that we thoroughly understand the relations of space and time implied by the Lorentz transformation, and therefore we shall consider this matter more deeply in this chapter.

The Lorentz transformation between the positions and times (x, y, z, t) as measured by an observer “standing still,” and the corresponding coordinates and time (x', y', z', t') measured inside a “moving” space ship, moving with velocity u are

$$\begin{aligned}x' &= \frac{x - ut}{\sqrt{1 - u^2/c^2}}, \\y' &= y, \\z' &= z, \\t' &= \frac{t - ux/c^2}{\sqrt{1 - u^2/c^2}}.\end{aligned}\tag{17.1}$$

Let us compare these equations with Eq. (11.5), which also relates measurements in two systems, one of which in this instance is *rotated* relative to the other:

$$\begin{aligned}x' &= x \cos \theta + y \sin \theta, \\y' &= y \cos \theta - x \sin \theta, \\z' &= z.\end{aligned}\tag{17.2}$$

In this particular case, Moe and Joe are measuring with axes having an angle θ between the x' - and x -axes. In each case, we note that the “primed” quantities are “mixtures” of the “unprimed” ones: the new x' is a mixture of x and y , and the new y' is also a mixture of x and y .

An analogy is useful: When we look at an object, there is an obvious thing we might call the “apparent width,” and another we might call the “depth.” But the two ideas, width and depth, are not *fundamental* properties of the object, because if we step aside and look at the same thing from a different angle, we get a different width and a different depth, and we may develop some formulas for computing the new ones from the old ones and the angles involved. Equations (17.2) are these formulas. One might say that a given depth is a kind of “mixture” of all depth and all width. If it were impossible ever to move, and we always saw a given object from the same position, then this whole business would be irrelevant—we would always see the “true” width and the “true” depth, and they would appear to have quite different qualities, because one appears as a subtended optical angle and the other involves some focusing of the eyes or even intuition; they would seem to be very different things and would never get mixed up. It is because we *can* walk around that we realize that depth and width are, somehow or other, just two different aspects of the same thing.

Can we not look at the Lorentz transformations in the same way? Here also we have a mixture—of positions and the time. A difference between a space measurement and a time measurement produces a new space measurement. In other words, in the space measurements of one man there is mixed in a little bit of the time, as seen by the other. Our analogy permits us to generate this idea: The “reality” of

17-1 The geometry of space-time

17-2 Space-time intervals

17-3 Past, present, and future

17-4 More about four-vectors

17-5 Four-vector algebra

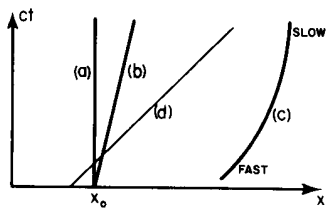


Fig. 17-1. Three particle paths in space-time: (a) a particle at rest at $x = x_0$; (b) a particle which starts at $x = x_0$ and moves with constant speed; (c) a particle which starts at high speed but slows down.

an object that we are looking at is somehow greater (speaking crudely and intuitively) than its “width” and its “depth” because *they* depend upon *how* we look at it; when we move to a new position, our brain immediately recalculates the width and the depth. But our brain does not immediately recalculate coordinates and time when we move at high speed, because we have had no effective experience of going nearly as fast as light to appreciate the fact that time and space are also of the same nature. It is as though we were always stuck in the position of having to look at just the width of something, not being able to move our heads appreciably one way or the other; if we could, we understand now, we would see some of the other man’s time—we would see “behind,” so to speak, a little bit.

Thus we shall try to think of objects in a new kind of world, of space and time mixed together, in the same sense that the objects in our ordinary space-world are real, and can be looked at from different directions. We shall then consider that objects occupying space and lasting for a certain length of time occupy a kind of a “blob” in a new kind of world, and that we look at this “blob” from different points of view when we are moving at different velocities. This new world, this geometrical entity in which the “blobs” exist by occupying position and taking up a certain amount of time, is called *space-time*. A given point (x, y, z, t) in space-time is called an *event*. Imagine, for example, that we plot the x -positions horizontally, y and z in two other directions, both mutually at “right angles” and at “right angles” to the paper (!), and time, vertically. Now, how does a moving particle, say, look on such a diagram? If the particle is standing still, then it has a certain x , and as time goes on, it has the same x , the same x , the same x ; so its “path” is a line that runs parallel to the t -axis (Fig. 17-1 a). On the other hand, if it drifts outward, then as the time goes on x increases (Fig. 17-1 b). So a particle, for example, which starts to drift out and then slows up should have a motion something like that shown in Fig. 17-1(c). A particle, in other words, which is permanent and does not disintegrate is represented by a line in space-time. A particle which disintegrates would be represented by a forked line, because it would turn into two other things which would start from that point.

What about light? Light travels at the speed c , and that would be represented by a line having a certain fixed slope (Fig. 17-1 d).

Now according to our new idea, if a given event occurs to a particle, say if it suddenly disintegrates at a certain space-time point into two new ones which follow some new tracks, and this interesting event occurred at a certain value of x and a certain value of t , then we would expect that, if this makes any sense, we just have to take a new pair of axes and turn them, and that will give us the new t and the new x in our new system, as shown in Fig. 17-2(a). But this is wrong, because Eq. (17.1) is not *exactly* the same mathematical transformation as Eq. (17.2). Note, for example, the difference in sign between the two, and the fact that one is written in terms of $\cos \theta$ and $\sin \theta$, while the other is written with algebraic quantities. (Of course, it is not impossible that the algebraic quantities could be written as cosine and sine, but actually they cannot.) But still, the two expressions *are* very similar. As we shall see, it is not really possible to think of space-time as a real, ordinary geometry because of that difference in sign. In fact, although we shall not emphasize this point, it turns out that a man who is moving has to use a set of axes which are inclined equally to the light ray, using a special kind of projection parallel to the x' - and t' -axes, for his x' and t' , as shown in Fig. 17-2(b). We shall not deal with the geometry, since it does not help much; it is easier to work with the equations.

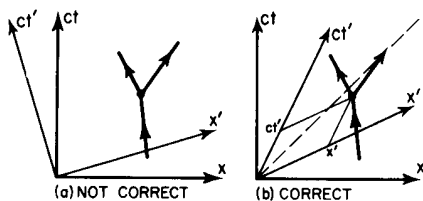


Fig. 17-2. Two views of a disintegrating particle.

17-2 Space-time intervals

Although the geometry of space-time is not Euclidean in the ordinary sense, there *is* a geometry which is very similar, but peculiar in certain respects. If this idea of geometry is right, there ought to be some functions of coordinates and time which are independent of the coordinate system. For example, under ordinary rotations, if we take two points, one at the origin, for simplicity, and the other one somewhere else, both systems would have the same origin, and the distance from

here to the other point is the same in both. That is one property that is independent of the particular way of measuring it. The square of the distance is $x^2 + y^2 + z^2$. Now what about space-time? It is not hard to demonstrate that we have here, also, something which stays the same, namely, the combination $c^2t^2 - x^2 - y^2 - z^2$ is the same before and after the transformation:

$$c^2t'^2 - x'^2 - y'^2 - z'^2 = c^2t^2 - x^2 - y^2 - z^2. \quad (17.3)$$

This quantity is therefore something which, like the distance, is “real” in some sense; it is called the *interval* between the two space-time points, one of which is, in this case, at the origin. (Actually, of course, it is the interval squared, just as $x^2 + y^2 + z^2$ is the distance squared.) We give it a different name because it is in a different geometry, but the interesting thing is only that some signs are reversed and there is a c in it.

Let us get rid of the c ; that is an absurdity if we are going to have a wonderful space with x 's and y 's that can be interchanged. One of the confusions that could be caused by someone with no experience would be to measure widths, say, by the angle subtended at the eye, and measure depth in a different way, like the strain on the muscles needed to focus them, so that the depths would be measured in feet and the widths in meters. Then one would get an enormously complicated mess of equations in making transformations such as (17.2), and would not be able to see the clarity and simplicity of the thing for a very simple technical reason, that the same thing is being measured in two different units. Now in Eqs. (17.1) and (17.3) nature is telling us that time and space are equivalent; time becomes space; *they should be measured in the same units*. What distance is a “second”? It is easy to figure out from (17.3) what it is. It is 3×10^8 meters, *the distance that light would go in one second*. In other words, if we were to measure all distances and times in the same units, seconds, then our unit of distance would be 3×10^8 meters, and the equations would be simpler. Or another way that we could make the units equal is to measure time in meters. What is a meter of time? A meter of time is the time it takes for light to go one meter, and is therefore $1/3 \times 10^{-8}$ sec, or 3.3 billionths of a second! We would like, in other words, to put all our equations in a system of units in which $c = 1$. If time and space are measured in the same units, as suggested, then the equations are obviously much simplified. They are

$$\begin{aligned} x' &= \frac{x - ut}{\sqrt{1 - u^2}}, \\ y' &= y, \\ z' &= z, \\ t' &= \frac{t - ux}{\sqrt{1 - u^2}}. \end{aligned} \quad (17.4)$$

$$t'^2 - x'^2 - y'^2 - z'^2 = t^2 - x^2 - y^2 - z^2. \quad (17.5)$$

If we are ever unsure or “frightened” that after we have this system with $c = 1$ we shall never be able to get our equations right again, the answer is quite the opposite. It is much easier to remember them without the c 's in them, and it is always easy to put the c 's back, by looking after the dimensions. For instance, in $\sqrt{1 - u^2}$, we know that we cannot subtract a velocity squared, which has units, from the pure number 1, so we know that we must divide u^2 by c^2 in order to make that unitless, and that is the way it goes.

The difference between space-time and ordinary space, and the character of an interval as related to the distance, is very interesting. According to formula (17.5), if we consider a point which in a given coordinate system had zero time, and only space, then the interval squared would be negative and we would have an imaginary interval, the square root of a negative number. Intervals can be either real or imaginary in the theory. The square of an interval may be either positive or negative, unlike distance, which has a positive square. When an interval is imaginary, we say that the two points have a *space-like interval* between them

(instead of imaginary), because the interval is more like space than like time. On the other hand, if two objects are at the same place in a given coordinate system, but differ only in time, then the square of the time is positive and the distances are zero and the interval squared is positive; this is called a *time-like interval*. In our diagram of space-time, therefore, we would have a representation something like this: at 45° there are two lines (actually, in four dimensions these will be “cones,” called light cones) and points on these lines are all at zero interval from the origin. Where light goes from a given point is always separated from it by a zero interval, as we see from Eq. (17.5). Incidentally, we have just proved that if light travels with speed c in one system, it travels with speed c in another, for if the interval is the same in both systems, i.e., zero in one and zero in the other, then to state that the propagation speed of light is invariant is the same as saying that the interval is zero.

17-3 Past, present, and future

The space-time region surrounding a given space-time point can be separated into three regions, as shown in Fig. 17-3. In one region we have space-like intervals, and in two regions, time-like intervals. Physically, these three regions into which space-time around a given point is divided have an interesting physical relationship to that point: a physical object or a signal can get from a point in region 2 to the event O by moving along at a speed less than the speed of light. Therefore events in this region can affect the point O , can have an influence on it from the past. In fact, of course, an object at P on the negative t -axis is precisely in the “past” with respect to O ; it is the same space-point as O , only earlier. What happened there then, affects O now. (Unfortunately, that is the way life is.) Another object at Q can get to O by moving with a certain speed less than c , so if this object were in a space ship and moving, it would be, again, the past of the same space-point. That is, in another coordinate system, the axis of time might go through both O and Q . So all points of region 2 are in the “past” of O , and anything that happens in this region *can* affect O . Therefore region 2 is sometimes called the *affective past*, or affecting past; it is the locus of all events which can affect point O in any way.

Region 3, on the other hand, is a region which we can affect *from* O , we can “hit” things by shooting “bullets” out at speeds less than c . So this is the world whose future can be affected by us, and we may call that the *affective future*. Now the interesting thing about all the rest of space-time, i.e., region 1, is that we can neither affect it now *from* O , nor can it affect us now *at* O , because nothing can go faster than the speed of light. Of course, what happens at R *can* affect us *later*; that is, if the sun is exploding “right now,” it takes eight minutes before we know about it, and it cannot possibly affect us before then.

What we mean by “right now” is a mysterious thing which we cannot define and we cannot affect, but it can affect us later, or we could have affected it if we had done something far enough in the past. When we look at the star Alpha Centauri, we see it as it was four years ago; we might wonder what it is like “now.” “Now” means at the same time from our special coordinate system. We can only see Alpha Centauri by the light that has come from our past, up to four years ago, but we do not know what it is doing “now”; it will take four years before what it is doing “now” can affect us. Alpha Centauri “now” is an idea or concept of our mind; it is not something that is really definable physically at the moment, because we have to wait to observe it; we cannot even define it right “now.” Furthermore, the “now” depends on the coordinate system. If, for example, Alpha Centauri were moving, an observer there would not agree with us because he would put his axes at an angle, and his “now” would be a *different* time. We have already talked about the fact that simultaneity is not a unique thing.

There are fortune tellers, or people who tell us they can know the future, and there are many wonderful stories about the man who suddenly discovers that he has knowledge about the affective future. Well, there are lots of paradoxes produced by that because if we know something is going to happen, then we can make

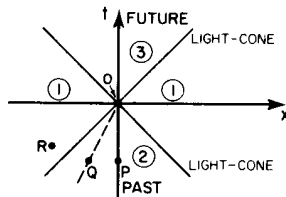


Fig. 17-3. The space-time region surrounding a point at the origin.

sure we will avoid it by doing the right thing at the right time, and so on. But actually there is no fortune teller who can even tell us the *present!* There is no one who can tell us what is really happening right now, at any reasonable distance, because that is unobservable. We might ask ourselves this question, which we leave to the student to try to answer: Would any paradox be produced if it were suddenly to become possible to know things that are in the space-like intervals of region 1?

17-4 More about four-vectors

Let us now return to our consideration of the analogy of the Lorentz transformation and rotations of the space axes. We have learned the utility of collecting together other quantities which have the same transformation properties as the coordinates, to form what we call *vectors*, directed lines. In the case of ordinary rotations, there are many quantities that transform the same way as x , y , and z under rotation: for example, the velocity has three components, an x , y , and z -component; when seen in a different coordinate system, none of the components is the same, instead they are all transformed to new values. But, somehow or other, the velocity "itself" has a greater reality than do any of its particular components, and we represent it by a directed line.

We therefore ask: Is it or is it not true that there are quantities which transform, or which are related, in a moving system and in a nonmoving system, in the same way as x , y , z , and t ? From our experience with vectors, we know that three of the quantities, like x , y , z , would constitute the three components of an ordinary space-vector, but the fourth quantity would look like an ordinary scalar under space rotation, because it does not change so long as we do not go into a moving coordinate system. Is it possible, then, to associate with some of our known "three-vectors" a fourth object, that we could call the "time component," in such a manner that the four objects together would "rotate" the same way as position and time in space-time? We shall now show that there is, indeed, at least one such thing (there are many of them, in fact): *the three components of momentum, and the energy as the time component, transform together* to make what we call a "four-vector." In demonstrating this, since it is quite inconvenient to have to write c 's everywhere, we shall use the same trick concerning units of the energy, the mass, and the momentum, that we used in Eq. (17.4). Energy and mass, for example, differ only by a factor c^2 which is merely a question of units, so we can say energy is the mass. Instead of having to write the c^2 , we put $E = m$, and then, of course, if there were any trouble we would put in the right amounts of c so that the units would straighten out in the last equation, but not in the intermediate ones.

Thus our equations for energy and momentum are

$$\begin{aligned} E &= m = m_0/\sqrt{1-v^2}, \\ \mathbf{p} &= m\mathbf{v} = m_0\mathbf{v}/\sqrt{1-v^2}. \end{aligned} \quad (17.6)$$

Also in these units, we have

$$E^2 - p^2 = m_0^2. \quad (17.7)$$

For example, if we measure energy in electron volts, what does a mass of 1 electron volt mean? It means the mass whose rest energy is 1 electron volt, that is, m_0c^2 is one electron volt. For example, the rest mass of an electron is 0.511×10^6 ev.

Now what would the momentum and energy look like in a new coordinate system? To find out, we shall have to transform Eq. (17.6), which we can do because we know how the velocity transforms. Suppose that, as we measure it, an object has a velocity v , but we look upon the same object from the point of view of a space ship which itself is moving with a velocity u , and in that system we use a prime to designate the corresponding thing. In order to simplify things at first, we shall take the case that the velocity v is in the direction of u . (Later, we can do the more general case.) What is v' , the velocity as seen from the space ship? It is the

composite velocity, the “difference” between v and u . By the law which we worked out before,

$$v' = \frac{v - u}{1 - uv}. \quad (17.8)$$

Now let us calculate the new energy E' , the energy as the fellow in the space ship would see it. He would use the same rest mass, of course, but he would use v' for the velocity. What we have to do is square v' , subtract it from one, take the square root, and take the reciprocal:

$$\begin{aligned} v'^2 &= \frac{v^2 - 2uv + u^2}{1 - 2uv + u^2v^2}, \\ 1 - v'^2 &= \frac{1 - 2uv + u^2v^2 - v^2 + 2uv - u^2}{1 - 2uv + u^2v^2} \\ &= \frac{1 - v^2 - u^2 + u^2v^2}{1 - 2uv + u^2v^2} \\ &= \frac{(1 - v^2)(1 - u^2)}{(1 - uv)^2}. \end{aligned}$$

Therefore

$$\frac{1}{\sqrt{1 - v'^2}} = \frac{1 - uv}{\sqrt{1 - v^2} \sqrt{1 - u^2}}. \quad (17.9)$$

The energy E' is then simply m_0 times the above expression. But we want to express the energy in terms of the unprimed energy and momentum, and we note that

$$E' = \frac{m_0 - m_0uv}{\sqrt{1 - v^2} \sqrt{1 - u^2}} = \frac{(m_0/\sqrt{1 - v^2}) - (m_0v/\sqrt{1 - v^2})u}{\sqrt{1 - u^2}},$$

or

$$E' = \frac{E - up_x}{\sqrt{1 - u^2}}, \quad (17.10)$$

which we recognize as being exactly of the same form as

$$t' = \frac{t - ux}{\sqrt{1 - u^2}}.$$

Next we must find the new momentum p'_x . This is just the energy E times v' , and is also simply expressed in terms of E and p :

$$p'_x = E'v' = \frac{m_0(1 - uv)}{\sqrt{1 - v^2} \sqrt{1 - u^2}} \cdot \frac{v - u}{(1 - uv)} = \frac{m_0v - m_0u}{\sqrt{1 - v^2} \sqrt{1 - u^2}}.$$

Thus

$$p'_x = \frac{p_x - uE}{\sqrt{1 - u^2}}, \quad (17.11)$$

which we recognize as being of precisely the same form as

$$x' = \frac{x - ut}{\sqrt{1 - u^2}}.$$

Thus the transformations for the new energy and momentum in terms of the old energy and momentum are exactly the same as the transformations for t' in terms of t and x , and x' in terms of x and t : all we have to do is, every time we see t in (17.4) substitute E , and every time we see x substitute p_x , and then the equations (17.4) will become the same as Eqs. (17.10) and (17.11). This would imply, if everything works right, an additional rule that $p'_y = p_y$ and that $p'_z = p_z$. To prove this would require our going back and studying the case of motion up and down. Actually, we did study the case of motion up and down in the last

chapter. We analyzed a complicated collision and we noticed that, in fact, the transverse momentum is *not* changed when viewed from a moving system; so we have already verified that $p'_y = p_y$ and $p'_z = p_z$. The complete transformation, then, is

$$\begin{aligned} p'_x &= \frac{p_x - uE}{\sqrt{1 - u^2}}, \\ p'_y &= p_y, \\ p'_z &= p_z \\ E' &= \frac{E - up_x}{\sqrt{1 - u^2}}. \end{aligned} \quad (17.12)$$

In these transformations, therefore, we have discovered four quantities which transform like x , y , z , and t , and which we call the *four-vector momentum*. Since the momentum is a four-vector, it can be represented on a space-time diagram of a moving particle as an "arrow" tangent to the path, as shown in Fig. 17-4. This arrow has a time component equal to the energy, and its space components represent its three-vector momentum; this arrow is more "real" than either the energy or the momentum, because those just depend on how we look at the diagram.

17-5 Four-vector algebra

The notation for four-vectors is different than it is for three-vectors. In the case of three-vectors, if we were to talk about the ordinary three-vector momentum we would write it \mathbf{p} . If we wanted to be more specific, we could say it has three components which are, for the axes in question, p_x , p_y , and p_z , or we could simply refer to a general component as p_i , and say that i could either be x , y , or z , and that these are the three components; that is, imagine that i is any one of three directions, x , y , or z . The notation that we use for four-vectors is analogous to this: we write p_μ for the four-vector, and μ stands for the *four* possible directions t , x , y , or z .

We could, of course, use any notation we want; do not laugh at notations; invent them, they are powerful. In fact, mathematics is, to a large extent, invention of better notations. The whole idea of a four-vector, in fact, is an improvement in notation so that the transformations can be remembered easily. A_μ , then, is a general four-vector, but for the special case of momentum, the p_t is identified as the energy, p_x is the momentum in the x -direction, p_y is that in the y -direction, and p_z is that in the z -direction. To add four-vectors, we add the corresponding components.

If there is an equation among four-vectors, then the equation is true for *each component*. For instance, if the law of conservation of three-vector momentum is to be true in particle collisions, i.e., if the sum of the momenta for a large number of interacting or colliding particles is to be a constant, that must mean that the sums of all momenta in the x -direction, in the y -direction, and in the z -direction, for all the particles, must each be constant. This law alone would be impossible in relativity because it is *incomplete*; it is like talking about only two of the components of a three-vector. It is incomplete because if we rotate the axes, we mix the various components, so we must include all three components in our law. Thus, in relativity, we must complete the law of conservation of momentum by extending it to include the *time* component. This is *absolutely necessary* to go with the other three, or there cannot be relativistic invariance. The *conservation of energy* is the fourth equation which goes with the conservation of momentum to make a valid four-vector relationship in the geometry of space and time. Thus the law of conservation of energy and momentum in four-dimensional notation is

$$\sum_{\text{particles in}} p_\mu = \sum_{\text{particles out}} p_\mu \quad (17.13)$$

or, in a slightly different notation

$$\sum_i p_{i\mu} = \sum_j p_{j\mu}, \quad (17.14)$$

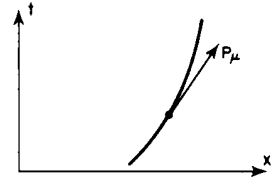


Fig. 17-4. The four-vector momentum of a particle.

where $i = 1, 2, \dots$ refers to the particles going into the collision, $j = 1, 2, \dots$ refers to the particles coming out of the collision, and $\mu = x, y, z$, or t . You say, "In which axes?" It makes no difference. The law is true for each component, using *any* axes.

In vector analysis we discussed one other thing, the dot product of two vectors. Let us now consider the corresponding thing in space-time. In ordinary rotation we discovered there was an unchanged quantity $x^2 + y^2 + z^2$. In four dimensions, we find that the corresponding quantity is $t^2 - x^2 - y^2 - z^2$ (Eq. 17.3). How can we write that? One way would be to write some kind of four-dimensional thing with a square dot between, like $A_\mu \diamond B_\mu$; one of the notations which is actually used is

$$\sum'_\mu A_\mu A_\mu = A_t^2 - A_x^2 - A_y^2 - A_z^2. \quad (17.15)$$

The prime on \sum means that the first term, the "time" term, is positive, but the other three terms have minus signs. This quantity, then, will be the same in any coordinate system, and we may call it the square of the length of the four-vector. For instance, what is the square of the length of the four-vector momentum of a single particle? This will be equal to $p_t^2 - p_x^2 - p_y^2 - p_z^2$ or, in other words, $E^2 - p^2$, because we know that p_t is E . What is $E^2 - p^2$? It must be something which is the same in every coordinate system. In particular, it must be the same for a coordinate system which is moving right along with the particle, in which the particle is standing still. If the particle is standing still, it would have no momentum. So in that coordinate system, it is purely its energy, which is the same as its rest mass. Thus $E^2 - p^2 = m_0^2$. So we see that the square of the length of this vector, the four-vector momentum, is equal to m_0^2 .

From the square of a vector, we can go on to invent the "dot product," or the product which is a scalar: if a_μ is one four-vector and b_μ is another four-vector, then the scalar product is

$$\sum'_\mu a_\mu b_\mu = a_t b_t - a_x b_x - a_y b_y - a_z b_z. \quad (17.16)$$

It is the same in all coordinate systems.

Finally, we shall mention certain things whose rest mass m_0 is zero. A photon of light, for example. A photon is like a particle, in that it carries an energy and a momentum. The energy of a photon is a certain constant, called Planck's constant, times the frequency of the photon: $E = h\nu$. Such a photon also carries a momentum, and the momentum of a photon (or of any other particle, in fact) is h divided by the wavelength: $p = h/\lambda$. But, for a photon, there is a definite relationship between the frequency and the wavelength: $\nu = c/\lambda$. (The number of waves per second, times the wavelength of each, is the distance that the light goes in one second, which, of course, is c .) Thus we see immediately that the energy of a photon must be the momentum times c , or if $c = 1$, *the energy and momentum are equal*. That is to say, the rest mass is zero. Let us look at that again; that is quite curious. If it is a particle of zero rest mass, what happens when it stops? *It never stops!* It always goes at the speed c . The usual formula for energy is $m_0/\sqrt{1 - v^2}$. Now can we say that $m_0 = 0$ and $v = 1$, so the energy is 0? We *cannot* say that it is zero; the photon really can (and does) have energy even though it has no rest mass, but this it possesses by perpetually going at the speed of light!

We also know that the momentum of any particle is equal to its total energy times its velocity: if $c = 1$, $p = vE$ or, in ordinary units, $p = vE/c^2$. For any particle moving at the speed of light, $p = E$ if $c = 1$. The formulas for the energy of a photon as seen from a moving system are, of course, given by Eq. (17.12), but for the momentum we must substitute the energy times c (or times 1 in this case). The different energies after transformation means that there are different frequencies. This is called the Doppler effect, and one can calculate it easily from Eq. (17.12), using also $E = p$ and $E = h\nu$.

As Minkowski said, "Space of itself, and time of itself will sink into mere shadows, and only a kind of union between them shall survive."

Rotation in Two Dimensions

18-1 The center of mass

In the previous chapters we have been studying the mechanics of points, or small particles whose internal structure does not concern us. For the next few chapters we shall study the application of Newton's laws to more complicated things. When the world becomes more complicated, it also becomes more interesting, and we shall find that the phenomena associated with the mechanics of a more complex object than just a point are really quite striking. Of course these phenomena involve nothing but combinations of Newton's laws, but it is sometimes hard to believe that only $F = ma$ is at work.

The more complicated objects we deal with can be of several kinds: water flowing, galaxies whirling, and so on. The simplest "complicated" object to analyze, at the start, is what we call a *rigid body*, a solid object that is turning as it moves about. However, even such a simple object may have a most complex motion, and we shall therefore first consider the simplest aspects of such motion, in which an extended body rotates about a *fixed axis*. A given point on such a body then moves in a plane perpendicular to this axis. Such rotation of a body about a fixed axis is called *plane rotation* or rotation in two dimensions. We shall later generalize the results to three dimensions, but in doing so we shall find that, unlike the case of ordinary particle mechanics, rotations are subtle and hard to understand unless we first get a solid grounding in two dimensions.

The first interesting theorem concerning the motion of complicated objects can be observed at work if we throw an object made of a lot of blocks and spokes, held together by strings, into the air. Of course we know it goes in a parabola, because we studied that for a particle. But now our object is *not* a particle; it wobbles and it jiggles, and so on. It does go in a parabola though; one can see that. *What* goes in a parabola? Certainly not the point on the corner of the block, because that is jiggling about; neither is it the end of the wooden stick, or the middle of the wooden stick, or the middle of the block. But *something* goes in a parabola, there is an effective "center" which moves in a parabola. So our first theorem about complicated objects is to demonstrate that there *is* a mean position which is mathematically definable, but not necessarily a point of the material itself, which goes in a parabola. That is called the theorem of the center of the mass, and the proof of it is as follows.

We may consider any object as being made of lots of little particles, the atoms, with various forces among them. Let i represent an index which defines one of the particles. (There are millions of them, so i goes to 10^{23} , or something.) Then the force on the i th particle is, of course, the mass times the acceleration of that particle:

$$\mathbf{F}_i = m_i(d^2\mathbf{r}_i/dt^2). \quad (18.1)$$

In the next few chapters our moving objects will be ones in which all the parts are moving at speeds very much slower than the speed of light, and we shall use the nonrelativistic approximation for all quantities. In these circumstances the mass is constant, so that

$$\mathbf{F}_i = d^2(m_i\mathbf{r}_i)/dt^2 \quad (18.2)$$

If we now add the force on all the particles, that is, if we take the sum of all the \mathbf{F}_i 's for all the different indexes, we get the total force, \mathbf{F} . On the other side of the

18-1 The center of mass

18-2 Rotation of a rigid body

18-3 Angular momentum

18-4 Conservation of angular momentum

equation, we get the same thing as though we added before the differentiation:

$$\sum_i \mathbf{F}_i = \mathbf{F} = \frac{d^2(\sum_i m_i \mathbf{r}_i)}{dt^2}. \quad (18.3)$$

Therefore the total force is the second derivative of the masses times their positions, added together.

Now the total force on all the particles is the same as the *external* force. Why? Although there are all kinds of forces on the particles because of the strings, the wiggles, the pullings and pushings, and the atomic forces, and who knows what, and we have to add all these together, we are rescued by Newton's Third Law. Between any two particles the action and reaction are equal, so that when we add all the equations together, if any two particles have forces between them it cancels out in the sum; therefore the net result is only those forces which arise from other particles which are not included in whatever object we decide to sum over. So if Eq. (18.3) is the sum over a certain number of the particles, which together are called "the object," then the *external* force on the total object is equal to the sum of *all* the forces on all its constituent particles.

Now it would be nice if we could write Eq. (18.3) as the total mass times some acceleration. We can. Let us say M is the sum of all the masses, i.e., the total mass. Then if we *define* a certain vector \mathbf{R} to be

$$\mathbf{R} = \sum_i m_i \mathbf{r}_i / M, \quad (18.4)$$

Eq. (18.3) will be simply

$$\mathbf{F} = d^2(M\mathbf{R})/dt^2 = M(d^2\mathbf{R}/dt^2), \quad (18.5)$$

since M is a constant. Thus we find that the external force is the total mass times the acceleration of an imaginary point whose location is \mathbf{R} . This point is called the *center of mass* of the body. It is a point somewhere in the "middle" of the object, a kind of average \mathbf{r} in which the different \mathbf{r}_i 's have weights or importances proportional to the masses.

We shall discuss this important theorem in more detail in a later chapter, and we shall therefore limit our remarks to two points: First, if the external forces are zero, if the object were floating in empty space, it might whirl, and jiggle, and twist, and do all kinds of things. But the *center of mass*, this artificially invented, calculated position, somewhere in the middle, *will move with a constant velocity*. In particular, if it is initially at rest, it will stay at rest. So if we have some kind of a box, perhaps a space ship, with people in it, and we calculate the location of the center of mass and find it is standing still, then the center of mass will continue to stand still if no external forces are acting on the box. Of course, the space ship may move a little in space, but that is because the people are walking back and forth inside; when one walks toward the front, the ship goes toward the back so as to keep the average position of all the masses in exactly the same place.

Is rocket propulsion therefore absolutely impossible because one cannot move the center of mass? No; but of course we find that to propel an interesting part of the rocket, an uninteresting part must be thrown away. In other words, if we start with a rocket at zero velocity and we spit some gas out the back end, then this little blob of gas goes one way as the rocket ship goes the other, but the center of mass is still exactly where it was before. So we simply move the part that we are interested in against the part we are not interested in.

The second point concerning the center of mass, which is the reason we introduced it into our discussion at this time, is that it may be treated separately from the "internal" motions of an object, and may therefore be ignored in our discussion of rotation.

18-2 Rotation of a rigid body

Now let us discuss rotations. Of course an ordinary object does not simply rotate, it wobbles, shakes, and bends, so to simplify matters we shall discuss the motion of a nonexistent ideal object which we call a *rigid body*. This means an

object in which the forces between the atoms are so strong, and of such character, that the little forces that are needed to move it do not bend it. Its shape stays essentially the same as it moves about. If we wish to study the motion of such a body, and agree to ignore the motion of its center of mass, there is only one thing left for it to do, and that is to *turn*. We have to describe that. How? Suppose there is some line in the body which stays put (perhaps it includes the center of mass and perhaps not), and the body is rotating about this particular line as an axis. How do we define the rotation? That is easy enough, for if we mark a point somewhere on the object, anywhere except on the axis, we can always tell exactly where the object is, if we only know where this point has gone to. The only thing needed to describe the position of that point is an *angle*. So rotation consists of a study of the variations of the angle with time.

In order to study rotation, we observe the angle through which a body has turned. Of course, we are not referring to any particular angle *inside* the object itself; it is not that we draw some angle *on* the object. We are talking about the *angular change of the position* of the whole thing, from one time to another.

First, let us study the kinematics of rotations. The angle will change with time, and just as we talked about position and velocity in one dimension, we may talk about angular position and angular velocity in plane rotation. In fact, there is a very interesting relationship between rotation in two dimensions and one-dimensional displacement, in which almost every quantity has its analog. First, we have the angle θ which defines how far the body has gone *around*; this replaces the distance y , which defines how far it has gone *along*. In the same manner, we have a velocity of turning, $\omega = d\theta/dt$, which tells us how much the angle changes in a second, just as $v = ds/dt$ describes how fast a thing moves, or how far it moves in a second. If the angle is measured in radians, then the angular velocity ω will be so and so many radians per second. The greater the angular velocity, the faster the object is turning, the faster the angle changes. We can go on: we can differentiate the angular velocity with respect to time, and we can call $\alpha = d\omega/dt = d^2\theta/dt^2$ the angular acceleration. That would be the analog of the ordinary acceleration.

Now of course we shall have to relate the dynamics of rotation to the laws of dynamics of the particles of which the object is made, so we must find out how a particular particle moves when the angular velocity is such and such. To do this, let us take a certain particle which is located at a distance r from the axis and say it is in a certain location $P(x, y)$ at a given instant, in the usual manner (Fig. 18-1). If at a moment Δt later the angle of the whole object has turned through $\Delta\theta$, then this particle is carried with it. It is at the same radius away from O as it was before, but is carried to Q . The first thing we would like to know is how much the distance x changes and how much the distance y changes. If OP is called r , then the length PQ is $r\Delta\theta$, because of the way angles are defined. The change in x , then, is simply the projection of $r\Delta\theta$ in the x -direction:

$$\Delta x = -PQ \sin \theta = -r \Delta\theta \cdot (y/r) = -y \Delta\theta. \quad (18.6)$$

Similarly,

$$\Delta y = +x \Delta\theta. \quad (18.7)$$

If the object is turning with a given angular velocity ω , we find, by dividing both sides of (18.6) and (18.7) by Δt , that the velocity of the particle is

$$v_x = -\omega y \quad \text{and} \quad v_y = +\omega x. \quad (18.8)$$

Of course if we want to find the magnitude of the velocity, we just write

$$v = \sqrt{v_x^2 + v_y^2} = \sqrt{\omega^2 y^2 + \omega^2 x^2} = \omega \sqrt{x^2 + y^2} = \omega r. \quad (18.9)$$

It should not be mysterious that the value of the magnitude of this velocity is ωr ; in fact, it should be self-evident, because the distance that it moves is $r\Delta\theta$ and the distance it moves per second is $r\Delta\theta/\Delta t$, or $r\omega$.

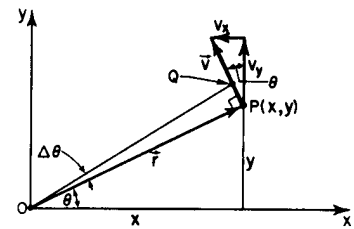


Fig. 18-1. Kinematics of two-dimensional rotation.

Let us now move on to consider the *dynamics* of rotation. Here a new concept, *force*, must be introduced. Let us inquire whether we can invent something which we shall call the *torque* (L. *torquere*, to twist) which bears the same relationship to rotation as force does to linear movement. A force is the thing that is needed to make linear motion, and the thing that makes something rotate is a “rotary force” or a “twisting force,” i.e., a torque. Qualitatively, a torque is a “twist”; what is a torque quantitatively? We shall get to the theory of torques quantitatively by studying the *work* done in turning an object, for one very nice way of defining a force is to say how much work it does when it acts through a given displacement. We are going to try to maintain the analogy between linear and angular quantities by equating the work that we do when we turn something a little bit when there are forces acting on it, to the *torque* times the *angle* it turns through. In other words, the definition of the torque is going to be so arranged that the theorem of work has an absolute analog: force times distance is work, and torque times angle is going to be work. That tells us what torque is. Consider, for instance, a rigid body of some kind with various forces acting on it, and an axis about which the body rotates. Let us at first concentrate on one force and suppose that this force is applied at a certain point (x, y) . How much work would be done if we were to turn the object through a very small angle? That is easy. The work done is

$$\Delta W = F_x \Delta x + F_y \Delta y. \quad (18.10)$$

We need only to substitute Eqs. (18.6) and (18.7) for Δx and Δy to obtain

$$\Delta W = (xF_y - yF_x) \Delta \theta. \quad (18.11)$$

That is, the amount of work that we have done is, in fact, equal to the angle through which we have turned the object, multiplied by a strange-looking combination of the force and the distance. This “strange combination” is what we call the torque. So, defining the change in work as the torque times the angle, we now have the formula for torque in terms of the forces. (Obviously, torque is not a completely new idea independent of Newtonian mechanics—torque must have a definite definition in terms of the force.)

When there are several forces acting, the work that is done is, of course, the sum of the works done by all the forces, so that ΔW will be a whole lot of terms, all added together, for all the forces, *each of which is proportional, however, to $\Delta \theta$* . We can take the $\Delta \theta$ outside and therefore can say that the change in the work is equal to the sum of all the torques due to all the different forces that are acting, times $\Delta \theta$. This sum we might call the total torque, τ . Thus torques add by the ordinary laws of algebra, but we shall later see that this is only because we are working in a plane. It is like one-dimensional kinematics, where the forces simply add algebraically, but only because they are all in the same direction. It is more complicated in three dimensions. Thus, for two-dimensional rotation,

$$\tau_i = x_i F_{yi} - y_i F_{xi} \quad (18.12)$$

and

$$\tau = \sum \tau_i. \quad (18.13)$$

It must be emphasized that the torque is about a given axis. If a different axis is chosen, so that all the x_i and y_i are changed, the value of the torque is (usually) changed too.

Now we pause briefly to note that our foregoing introduction of torque, through the idea of work, gives us a most important result for an object in equilibrium: if all the forces on an object are in balance both for translation and rotation, not only is the net *force* zero, but the total of all the *torques* is also zero, because if an object is in equilibrium, *no work is done by the forces for a small displacement*. Therefore, since $\Delta W = \tau \Delta \theta = 0$, the sum of all the torques must be zero. So there are two conditions for equilibrium: that the sum of the forces is zero, and that the sum of the torques is zero. Prove that it suffices to be sure that the sum of torques about any *one* axis (in two dimensions) is zero.

Now let us consider a single force, and try to figure out, geometrically, what this strange thing $x\bar{F}_y - y\bar{F}_x$ amounts to. In Fig. 18-2 we see a force \bar{F} acting at a point r . When the object has rotated through a small angle $\Delta\theta$, the work done, of course, is the component of force in the direction of the displacement times the displacement. In other words, it is only the tangential component of the force that counts, and this must be multiplied by the distance $r\Delta\theta$. Therefore we see that the torque is also equal to the tangential component of force (perpendicular to the radius) times the radius. That makes sense in terms of our ordinary idea of the torque, because if the force were completely radial, it would not put any "twist" on the body; it is evident that the twisting effect should involve only the part of the force which is not pulling out from the center, and that means the tangential component. Furthermore, it is clear that a given force is more effective on a long arm than near the axis. In fact, if we take the case where we push right on the axis, we are not twisting at all! So it makes sense that the amount of twist, or torque, is proportional both to the radial distance and to the tangential component of the force.

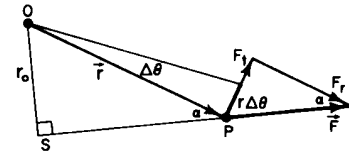


Fig. 18-2. The torque produced by a force.

There is still a third formula for the torque which is very interesting. We have just seen that the torque is the force times the radius times the sine of the angle α , in Fig. 18-2. But if we extend the line of action of the force and draw the line OS , the perpendicular distance to the line of action of the force (the *lever arm* of the force) we notice that this lever arm is shorter than r in just the same proportion as the tangential part of the force is less than the total force. Therefore the formula for the torque can also be written as the magnitude of the force times the length of the lever arm.

The torque is also often called the *moment* of the force. The origin of this term is obscure, but it may be related to the fact that "moment" is derived from the Latin *movimentum*, and that the capability of a force to move an object (using the force on a lever or crowbar) increases with the length of the lever arm. In mathematics "moment" means weighted by how far away it is from an axis.

18-3 Angular momentum

Although we have so far considered only the special case of a rigid body, the properties of torques and their mathematical relationships are interesting also even when an object is not rigid. In fact, we can prove a very remarkable theorem: just as external force is the rate of change of a quantity p , which we call the total momentum of a collection of particles, so the external torque is the rate of change of a quantity L which we call the *angular momentum* of the group of particles.

To prove this, we shall suppose that there is a system of particles on which there are some forces acting and find out what happens to the system as a result of the torques due to these forces. First, of course, we should consider just *one* particle. In Fig. 18-3 is one particle of mass m , and an axis O ; the particle is not necessarily rotating in a circle about O , it may be moving in an ellipse, like a planet going around the sun, or in some other curve. It is moving somehow, and there are forces on it, and it accelerates according to the usual formula that the x -component of force is the mass times the x -component of acceleration, etc. But let us see what the *torque* does. The torque equals $x\bar{F}_y - y\bar{F}_x$, and the force in the x - or y -direction is the mass times the acceleration in the x - or y -direction:

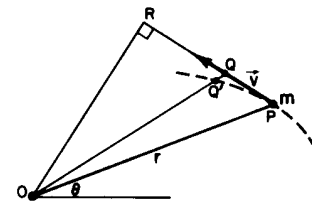


Fig. 18-3. A particle moves about an axis O .

$$\begin{aligned}\tau &= x\bar{F}_y - y\bar{F}_x \\ &= xm(d^2y/dt^2) - ym(d^2x/dt^2).\end{aligned}\quad (18.14)$$

Now, although this does not appear to be the derivative of any simple quantity, it is in fact the derivative of the quantity $xm(dy/dt) - ym(dx/dt)$:

$$\begin{aligned}\frac{d}{dt} \left[xm \left(\frac{dy}{dt} \right) - ym \left(\frac{dx}{dt} \right) \right] &= xm \left(\frac{d^2y}{dt^2} \right) + \left(\frac{dx}{dt} \right) m \left(\frac{dy}{dt} \right) \\ &\quad - ym \left(\frac{d^2x}{dt^2} \right) - \left(\frac{dy}{dt} \right) m \left(\frac{dx}{dt} \right) = xm \left(\frac{d^2y}{dt^2} \right) - ym \left(\frac{d^2x}{dt^2} \right).\end{aligned}\quad (18.15)$$

So it is true that the torque is the rate of change of something with time! So we pay attention to the “something,” we give it a name: we call it L , the angular momentum:

$$\begin{aligned} L &= xm(dy/dt) - ym(dx/dt) \\ &= xp_y - yp_x. \end{aligned} \quad (18.16)$$

Although our present discussion is nonrelativistic, the second form for L given above is relativistically correct. So we have found that there is also a rotational analog for the momentum, and that this analog, the angular momentum, is given by an expression in terms of the components of linear momentum that is just like the formula for torque in terms of the force components! Thus, if we want to know the angular momentum of a particle about an axis, we take only the component of the momentum that is tangential, and multiply it by the radius. In other words, what counts for angular momentum is not how fast it is going *away* from the origin, but how much it is going *around* the origin. Only the tangential part of the momentum counts for angular momentum. Furthermore, the farther out the line of the momentum extends, the greater the angular momentum. And also, because the geometrical facts are the same whether the quantity is labeled p or F , it is true that there is a lever arm (*not* the same as the lever arm of the force on the particle!) which is obtained by extending the line of the *momentum* and finding the perpendicular distance to the axis. Thus the angular momentum is the magnitude of the momentum times the momentum lever arm. So we have three formulas for angular momentum, just as we have three formulas for the torque:

$$\begin{aligned} L &= xp_y - yp_x \\ &= rp_{\text{tang}} \\ &= p \cdot \text{lever arm}. \end{aligned} \quad (18.17)$$

Like torque, angular momentum depends upon the position of the axis about which it is to be calculated.

Before proceeding to a treatment of more than one particle, let us apply the above results to a planet going around the sun. In which direction is the force? The force is toward the sun. What, then, is the torque on the object? Of course, this depends upon where we take the axis, but we get a very simple result if we take it at the sun itself, for the torque is the force times the lever arm, or the component of force perpendicular to r , times r . But there is no tangential force, so there is no torque about an axis at the sun! Therefore, the angular momentum of the planet going around the sun must remain constant. Let us see what that means. The tangential component of velocity, times the mass, times the radius, will be constant, because that is the angular momentum, and the rate of change of the angular momentum is the torque, and, in this problem, the torque is zero. Of course since the mass is also a constant, this means that the tangential velocity times the radius is a constant. But this is something we already knew for the motion of a planet. Suppose we consider a small amount of time Δt . How far will the planet move when it moves from P to Q (Fig. 18-3)? How much *area* will it sweep through? Disregarding the very tiny area $QQ'P$ compared with the much larger area OPQ , it is simply half the base PQ times the height, OR . In other words, the area that is swept through in unit time will be equal to the velocity times the lever arm of the velocity (times one-half). Thus the rate of change of area is proportional to the angular momentum, which is constant. So Kepler's law about equal areas in equal times is a word description of the statement of the law of conservation of angular momentum, when there is no torque produced by the force.

18-4 Conservation of angular momentum

Now we shall go on to consider what happens when there is a large number of particles, when an object is made of many pieces with many forces acting between them and on them from the outside. Of course, we already know that, about any given fixed axis, the torque on the i th particle (which is the force on the i th particle

times the lever arm of that force) is equal to the rate of change of the angular momentum of that particle, and that the angular momentum of the i th particle is its momentum times its momentum lever arm. Now suppose we add the torques τ_i for all the particles and call it the total torque τ . Then this will be the rate of change of the sum of the angular momenta of all the particles L_i , and that defines a new quantity which we call the total angular momentum L . Just as the total *momentum* of an object is the sum of the momenta of all the parts, so the angular momentum is the sum of the angular momenta of all the parts. Then the rate of change of the total L is the total torque:

$$\tau = \sum \tau_i = \sum \frac{dL_i}{dt} = \frac{dL}{dt}. \quad (18.18)$$

Now it might seem that the total torque is a complicated thing. There are all those internal forces and all the outside forces to be considered. But, if we take Newton's law of action and reaction to say, not simply that the action and reaction are equal, but also that they are *directed exactly oppositely along the same line* (Newton may or may not actually have said this, but he tacitly assumed it), then the two *torques* on the reacting objects, due to their mutual interaction, will be equal and opposite because the lever arms for any axis are equal. Therefore the internal torques balance out pair by pair, and so we have the remarkable theorem that the *rate of change of the total angular momentum about any axis is equal to the external torque about that axis!*

$$\tau = \sum \tau_i = \tau_{\text{ext}} = dL/dt. \quad (18.19)$$

Thus we have a very powerful theorem concerning the motion of large collections of particles, which permits us to study the over-all motion without having to look at the detailed machinery inside. This theorem is true for any collection of objects, whether they form a rigid body or not.

One extremely important case of the above theorem is the *law of conservation of angular momentum*: if no external torques act upon a system of particles, the angular momentum remains constant.

A special case of great importance is that of a rigid body, that is, an object of a definite shape that is just turning around. Consider an object that is fixed in its geometrical dimensions, and which is rotating about a fixed axis. Various parts of the object bear the same relationship to one another at all times. Now let us try to find the total angular momentum of this object. If the mass of one of its particles is m_i , and its position or location is at (x_i, y_i) , then the problem is to find the angular momentum of that particle, because the total angular momentum is the sum of the angular momenta of all such particles in the body. For an object going around in a circle, the angular momentum, of course, is the mass times the velocity times the distance from the axis, and the velocity is equal to the angular velocity times the distance from the axis:

$$L_i = m_i v_i r_i = m_i r_i^2 \omega, \quad (18.20)$$

or, summing over all the particles i , we get

$$L = I\omega, \quad (18.21)$$

where

$$I = \sum_i m_i r_i^2. \quad (18.22)$$

This is the analog of the law that the momentum is mass times velocity. Velocity is replaced by angular velocity, and we see that the mass is replaced by a new thing which we call the *moment of inertia* I , which is analogous to the mass. Equations (18.21) and (18.22) say that a body has inertia for turning which depends, not just on the masses, but on *how far away they are* from the axis. So, if we have two objects of the same mass, when we put the masses farther away from the axis, the inertia for turning will be higher. This is easily demonstrated by the apparatus

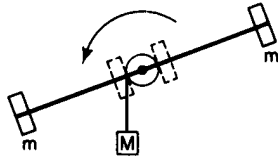


Fig. 18-4. The "inertia for turning" depends upon the lever arm of the masses.

shown in Fig. 18-4, where a weight M is kept from falling very fast because it has to turn the large weighted rod. At first, the masses m are close to the axis, and M speeds up at a certain rate. But when we change the moment of inertia by putting the two masses m much farther away from the axis, then we see that M accelerates much less rapidly than it did before, because the body has much more inertia against turning. The moment of inertia is the inertia against turning, and is the sum of the contributions of all the masses, times their distances *squared*, from the axis.

There is one important difference between mass and moment of inertia which is very dramatic. The mass of an object never changes, but its moment of inertia *can* be changed. If we stand on a frictionless rotatable stand with our arms outstretched, and hold some weights in our hands as we rotate slowly, we may change our moment of inertia by drawing our arms in, but our mass does not change. When we do this, all kinds of wonderful things happen, because of the law of the conservation of angular momentum: If the external torque is zero, then the angular momentum, the moment of inertia times omega, remains constant. Initially, we were rotating with a large moment of inertia I_1 at a low angular velocity ω_1 , and the angular momentum was $I_1\omega_1$. Then we changed our moment of inertia by pulling our arms in, say to a smaller value I_2 . Then the product $I\omega$, which has to stay the same because the total angular momentum has to stay the same, was $I_2\omega_2$. So $I_1\omega_1 = I_2\omega_2$. That is, if we *reduce* the moment of inertia, we have to *increase* the angular velocity.

Center of Mass; Moment of Inertia

19-1 Properties of the center of mass

In the previous chapter we found that if a great many forces are acting on a complicated mass of particles, whether the particles comprise a rigid or a nonrigid body, or a cloud of stars, or anything else, and we find the sum of all the forces (that is, of course, the external forces, because the internal forces balance out), then if we consider the body as a whole, and say it has a total mass M , there is a certain point "inside" the body, called the *center of mass*, such that the net resulting external force produces an acceleration of this point, just as though the whole mass were concentrated there. Let us now discuss the center of mass in a little more detail.

The location of the center of mass (abbreviated CM) is given by the equation

$$\mathbf{R}_{\text{CM}} = \frac{\sum m_i \mathbf{r}_i}{\sum m_i}. \quad (19.1)$$

This is, of course, a vector equation which is really three equations, one for each of the three directions. We shall consider only the x -direction, because if we can understand that one, we can understand the other two. What does $X_{\text{CM}} = \sum m_i x_i / \sum m_i$ mean? Suppose for a moment that the object is divided into little pieces, all of which have the same mass m ; then the total mass is simply the number N of pieces times the mass of one piece, say one gram, or any unit. Then this equation simply says that we add all the x 's, and then divide by the number of things that we have added: $X_{\text{CM}} = m \sum x_i / mN = \sum x_i / N$. In other words, X_{CM} is the average of all the x 's, if the masses are equal. But suppose one of them were twice as heavy as the others. Then in the sum, that x would come in twice. This is easy to understand, for we can think of this double mass as being split into two equal ones, just like the others; then in taking the average, of course, we have to count that x twice because there are two masses there. Thus X is the average position, in the x -direction, of all the masses, every mass being counted a number of times proportional to the mass, as though it were divided into "little grams." From this it is easy to prove that X must be somewhere between the largest and the smallest x , and, therefore lies inside the envelope including the entire body. It does not have to be in the *material* of the body, for the body could be a circle, like a hoop, and the center of mass is in the center of the hoop, not in the hoop itself.

Of course, if an object is symmetrical in some way, for instance, a rectangle, so that it has a plane of symmetry, the center of mass lies somewhere on the plane of symmetry. In the case of a rectangle there are two planes, and that locates it uniquely. But if it is just any symmetrical object, then the center of gravity lies somewhere on the axis of symmetry, because in those circumstances there are as many positive as negative x 's.

Another interesting proposition is the following very curious one. Suppose that we imagine an object to be made of two pieces, A and B (Fig. 19-1). Then the center of mass of the whole object can be calculated as follows. First, find the center of mass of piece A , and then of piece B . Also, find the total mass of each piece, M_A and M_B . Then consider a new problem, in which a *point* mass M_A is at the center of mass of object A , and another *point* mass M_B is at the center of mass of object B . The center of mass of these two point masses is then the center of mass of the whole object. In other words, if the centers of mass of various parts

19-1 Properties of the center of mass

19-2 Locating the center of mass

19-3 Finding the moment of inertia

19-4 Rotational kinetic energy

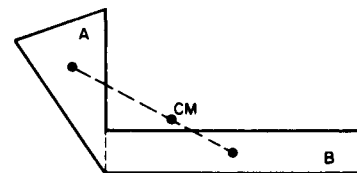


Fig. 19-1. The CM of a compound body lies on the line joining the CM's of the two composite parts.

of an object have been worked out, we do not have to start all over again to find the center of mass of the whole object; we just have to put the pieces together, treating each one as a point mass situated at the center of mass of that piece. Let us see why that is. Suppose that we wanted to calculate the center of mass of a complete object, some of whose particles are considered to be members of object A and some members of object B . The total sum $\sum m_i x_i$ can then be split into two pieces—the sum $\sum_A m_i x_i$ for the A object only, and the sum $\sum_B m_i x_i$ for object B only. Now if we were computing the center of mass of object A alone, we would have exactly the first of these sums, and we know that this by itself is $M_A X_A$, the total mass of all the particles in A times the position of the center of mass of A , because that is the theorem of the center of mass, applied to object A . In the same manner, just by looking at object B , we get $M_B X_B$, and of course, adding the two yields MX :

$$\begin{aligned} MX_{\text{CM}} &= \sum_A m_i x_i + \sum_B m_i x_i \\ &= M_A X_A + M_B X_B. \end{aligned} \quad (19.2)$$

Now since M is evidently the sum of M_A and M_B , we see that Eq. (19.2) can be interpreted as a special example of the center of mass formula for two point objects, one of mass M_A located at X_A and the other of mass M_B located at X_B .

The theorem concerning the motion of the center of mass is very interesting, and has played an important part in the development of our understanding of physics. Suppose we assume that Newton's law is right for the small component parts of a much larger object. Then this theorem shows that Newton's law is also correct for the larger object, even if we do not study the details of the object, but only the total force acting on it and its mass. In other words, Newton's law has the peculiar property that if it is right on a certain small scale, then it will be right on a larger scale. If we do not consider a baseball as a tremendously complex thing, made of myriads of interacting particles, but study only the motion of the center of mass and the external forces on the ball, we find $\mathbf{F} = m\mathbf{a}$, where \mathbf{F} is the external force on the baseball, m is its mass, and \mathbf{a} is the acceleration of its center of mass. So $\mathbf{F} = m\mathbf{a}$ is a law which reproduces itself on a larger scale. (There ought to be a good word, out of the Greek, perhaps, to describe a law which reproduces the same law on a larger scale.)

Of course, one might suspect that the first laws that would be discovered by human beings would be those that would reproduce themselves on a larger scale. Why? Because the actual scale of the fundamental gears and wheels of the universe are of atomic dimensions, which are so much finer than our observations that we are nowhere near that scale in our ordinary observations. So the first things that we would discover must be true for objects of no special size relative to an atomic scale. If the laws for small particles did not reproduce themselves on a larger scale, we would not discover those laws very easily. What about the reverse problem? Must the laws on a small scale be the same as those on a larger scale? Of course it is not necessarily so in nature, that at an atomic level the laws have to be the same as on a large scale. Suppose that the true laws of motion of atoms were given by some strange equation which does *not* have the property that when we go to a larger scale we reproduce the same law, but instead has the property that if we go to a larger scale, we can *approximate it by a certain expression* such that, if we extend that expression up and up, *it* keeps reproducing itself on a larger and larger scale. That is possible, and in fact that is the way it works. Newton's laws are the "tail end" of the atomic laws, extrapolated to a very large size. The actual laws of motion of particles on a fine scale are very peculiar, but if we take large numbers of them and compound them, they approximate, but *only* approximate, Newton's laws. Newton's laws then permit us to go on to a higher and higher scale, and it still seems to be the same law. In fact, it becomes more and more accurate as the scale gets larger and larger. This self-reproducing factor of Newton's laws is thus really not a fundamental feature of nature, but is an important historical feature. We would never discover the fundamental laws of the atomic particles at first observation because the first observations are much too crude. In fact, it turns

out that the fundamental atomic laws, which we call quantum mechanics, are quite different from Newton's laws, and are difficult to understand because all our direct experiences are with large-scale objects and the small-scale atoms behave like nothing we see on a large scale. So we cannot say, "An atom is just like a planet going around the sun," or anything like that. It is like *nothing* we are familiar with because there is *nothing like it*. As we apply quantum mechanics to larger and larger things, the laws about the behavior of many atoms together do *not* reproduce themselves, but produce *new laws*, which are Newton's laws, which then continue to reproduce themselves from, say, micro-microgram size, which still is billions and billions of atoms, on up to the size of the earth, and above.

Let us now return to the center of mass. The center of mass is sometimes called the center of gravity, for the reason that, in many cases, gravity may be considered uniform. Let us suppose that we have small enough dimensions that the gravitational force is not only proportional to the mass, but is everywhere parallel to some fixed line. Then consider an object in which there are gravitational forces on each of its constituent masses. Let m_i be the mass of one part. Then the gravitational force on that part is m_i times g . Now the question is, where can we apply a single force to balance the gravitational force on the whole thing, so that the entire object, if it is a rigid body, will not turn? The answer is that this force must go through the center of mass, and we show this in the following way. In order that the body will not turn, the torque produced by all the forces must add up to zero, because if there is a torque, there is a change of angular momentum, and thus a rotation. So we must calculate the total of all the torques on all the particles, and see how much torque there is about any given axis; it should be zero if this axis is at the center of mass. Now, measuring x horizontally and y vertically, we know that the torques are the forces in the y -direction, times the lever arm x (that is to say, the force times the lever arm around which we want to measure the torque). Now the total torque is the sum

$$\tau = \sum m_i g x_i = g \sum m_i x_i, \quad (19.3)$$

so if the total torque is to be zero, the sum $\sum m_i x_i$ must be zero. But $\sum m_i x_i = MX$, the total mass times the distance of the center of mass from the axis. Thus the x -distance of the center of mass from the axis is zero.

Of course, we have checked the result only for the x -distance, but if we use the true center of mass the object will balance in any position, because if we turned it 90 degrees, we would have y 's instead of x 's. In other words, when an object is supported at its center of mass, there is no torque on it because of a parallel gravitational field. In case the object is so large that the nonparallelism of the gravitational forces is significant, then the center where one must apply the balancing force is not simple to describe, and it departs slightly from the center of mass. That is why one must distinguish between the center of mass and the center of gravity. The fact that an object supported exactly at the center of mass will balance in all positions has another interesting consequence. If, instead of gravitation, we have a pseudoforce due to acceleration, we may use exactly the same mathematical procedure to find the position to support it so that there are no torques produced by the inertial force of acceleration. Suppose that the object is held in some manner inside a box, and that the box, and everything contained in it, is accelerating. We know that, from the point of view of someone at rest relative to this accelerating box, there will be an effective force due to inertia. That is, to make the object go along with the box, we have to push on it to accelerate it, and this force is "balanced" by the "force of inertia," which is a pseudoforce equal to the mass times the acceleration of the box. To the man in the box, this is the same situation as if the object were in a uniform gravitational field whose " g " value is equal to the acceleration a . Thus the inertial force due to accelerating an object has no torque about the center of mass.

This fact has a very interesting consequence. In an inertial frame that is not accelerating, the torque is always equal to the rate of change of the angular momentum. However, about an axis through the center of mass of an object which

is accelerating, it is *still true* that the torque is equal to the rate of change of the angular momentum. Even if the center of mass is accelerating, we may still choose one special axis, namely, one passing through the center of mass, such that it will still be true that the torque is equal to the rate of change of angular momentum around that axis. Thus the theorem that torque equals the rate of change of angular momentum is true in two general cases: (1) a fixed axis in inertial space, (2) an axis through the center of mass, even though the object may be accelerating.

19-2 Locating the center of mass

The mathematical techniques for the calculation of centers of mass are in the province of a mathematics course, and such problems provide good exercise in integral calculus. After one has learned calculus, however, and wants to know how to locate centers of mass, it is nice to know certain tricks which can be used to do so. One such trick makes use of what is called the theorem of Pappus. It works like this: if we take any closed area in a plane and generate a solid by moving it through space such that each point is always moved perpendicular to the plane of the area, the resulting solid has a total volume equal to the area of the cross section times the distance that the center of mass moved! Certainly this is true if we move the area in a straight line perpendicular to itself, but if we move it in a circle or in some other curve, then it generates a rather peculiar volume. For a curved path, the outside goes around farther, and the inside goes around less, and these effects balance out. So if we want to locate the center of mass of a plane sheet of uniform density, we can remember that the volume generated by spinning it about an axis is the distance that the center of mass goes around, times the area of the sheet.

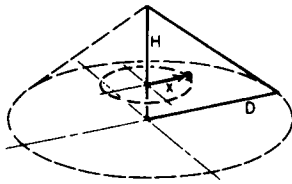


Fig. 19-2. A right triangle and a right circular cone generated by rotating the triangle.

For example, if we wish to find the center of mass of a right triangle of base D and height H (Fig. 19-2), we might solve the problem in the following way. Imagine an axis along H , and rotate the triangle about that axis through a full 360 degrees. This generates a cone. The distance that the x -coordinate of the center of mass has moved is $2\pi x$. The area which is being moved is the area of the triangle, $\frac{1}{2}HD$. So the x -distance of the center of mass times the area of the triangle is the volume swept out, which is of course $\pi D^2 H/3$. Thus $(2\pi x)(\frac{1}{2}HD) = 1/3\pi D^2 H$, or $x = D/3$. In a similar manner, by rotating about the other axis, or by symmetry, we find $y = H/3$. In fact, the center of mass of any uniform triangular area is where the three medians, the lines from the vertices through the centers of the opposite sides, all meet. That point is $1/3$ of the way along each median. *Clue:* Slice the triangle up into a lot of little pieces, each parallel to a base. Note that the median line bisects every piece, and therefore the center of mass must lie on this line.

Now let us try a more complicated figure. Suppose that it is desired to find the position of the center of mass of a uniform semicircular disc—a disc sliced in half. Where is the center of mass? For a full disc, it is at the center, of course, but a half-disc is more difficult. Let r be the radius and x be the distance of the center of mass from the straight edge of the disc. Spin it around this edge as axis to generate a sphere. Then the center of mass has gone around $2\pi x$, the area is $\pi r^2/2$ (because it is only half a circle). The volume generated is, of course, $4\pi r^3/3$, from which we find that

$$(2\pi x)(\frac{1}{2}\pi r^2) = 4\pi r^3/3,$$

or

$$x = 4r/3\pi.$$

There is another theorem of Pappus which is a special case of the above one, and therefore equally true. Suppose that, instead of the solid semicircular disc, we have a semicircular piece of wire with uniform mass density along the wire, and we want to find its center of mass. In this case there is no mass in the interior, only on the wire. Then it turns out that the *area* which is swept by a plane curved line, when it moves as before, is the distance that the center of mass moves times the *length* of the line. (The line can be thought of as a very narrow area, and the previous theorem can be applied to it.)

19-3 Finding the moment of inertia

Now let us discuss the problem of finding the *moments of inertia* of various objects. The formula for the moment of inertia about the z -axis of an object is

$$I = \sum m_i(x_i^2 + y_i^2)$$

or

$$I = \int (x^2 + y^2) dm = \int (x^2 + y^2) \rho dv. \quad (19.4)$$

That is, we must sum the masses, each one multiplied by the square of its distance ($x_i^2 + y_i^2$) from the axis. Note that it is not the three-dimensional distance, only the two-dimensional distance squared, even for a three-dimensional object. For the most part, we shall restrict ourselves to two-dimensional objects, but the formula for rotation about the z -axis is just the same in three dimensions.

As a simple example, consider a rod rotating about a perpendicular axis through one end (Fig. 19-3). Now we must sum all the masses times the x -distances squared (the y 's being all zero in this case). What we mean by "the sum," of course, is the integral of x^2 times the little elements of mass. If we divide the rod into small elements of length dx , the corresponding elements of mass are proportional to dx , and if dx were the length of the whole rod the mass would be M . Therefore

$$dm = M dx/L$$

and so

$$I = \int_0^L x^2 \frac{M dx}{L} = \frac{M}{L} \int_0^L x^2 dx = \frac{ML^2}{3}. \quad (19.5)$$

The dimensions of moment of inertia are always mass times length squared, so all we really had to work out was the factor $1/3$.

Now what is I if the rotation axis is at the center of the rod? We could just do the integral over again, letting x range from $-\frac{1}{2}L$ to $+\frac{1}{2}L$. But let us notice a few things about the moment of inertia. We can imagine the rod as two rods, each of mass $M/2$ and length $L/2$; the moments of inertia of the two small rods are equal, and are both given by the formula (19.5). Therefore the moment of inertia is

$$I = \frac{2(M/2)(L/2)^2}{3} = \frac{ML^2}{12}. \quad (19.6)$$

Thus it is much easier to turn a rod about its center, than to swing it around an end.

Of course, we could go on to compute the moments of inertia of various other bodies of interest. However, while such computations provide a certain amount of important exercise in the calculus, they are not basically of interest to us as such. There is, however, an interesting theorem which is very useful. Suppose we have an object, and we want to find its moment of inertia around some axis. That means we want the inertia needed to carry it by rotation about that axis. Now if we support the object on pivots at the center of mass, so that the object does not turn as it rotates about the axis (because there is no torque on it from inertial effects, and therefore it will not turn when we start moving it), then the forces needed to swing it around are the same as though all the mass were concentrated at the center of mass, and the moment of inertia would be simply $I_1 = MR_{CM}^2$, where R_{CM} is the distance from the axis to the center of mass. But of course that is not the right formula for the moment of inertia of an object which is really being rotated as it revolves, because not only is the center of it moving in a circle, which would contribute an amount I_1 to the moment of inertia, but also we must turn it about its center of mass. So it is not unreasonable that we must add to I_1 the moment of inertia I_c about the center of mass. So it is a good guess that the total moment of inertia about any axis will be

$$I = I_c + MR_{CM}^2. \quad (19.7)$$

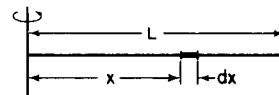


Fig. 19-3. A straight rod of length L rotating about an axis through one end.

This theorem is called the *parallel-axis theorem*, and may be easily proved. The moment of inertia about any axis is the mass times the sum of the x_i 's and the y_i 's, each squared: $I = \sum (x_i^2 + y_i^2)m_i$. We shall concentrate on the x 's, but of course the y 's work the same way. Now x is the distance of a particular point mass from the origin, but let us consider how it would look if we measured x' from the CM, instead of x from the origin. To get ready for this analysis, we write

$$x_i = x'_i + X_{\text{CM}}.$$

Then we just square this to find

$$x_i^2 = x'^2_i + 2X_{\text{CM}}x'_i + X_{\text{CM}}^2.$$

So, when this is multiplied by m_i and summed over all i , what happens? Taking the constants outside the summation sign, we get

$$I_x = \sum m_i x'^2_i + 2X_{\text{CM}} \sum m_i x'_i + X_{\text{CM}}^2 \sum m_i.$$

The third sum is easy; it is just MX_{CM}^2 . In the second sum there are two pieces, one of them is $\sum m_i x'_i$, which is the total mass times the x' -coordinate of the center of mass. But this contributes nothing, because x' is *measured* from the center of mass, and in these axes the average position of all the particles, weighted by the masses, is zero. The first sum, of course, is the x part of I_c . Thus we arrive at Eq. (19.7), just as we guessed.

Let us check (19.7) for one example. Let us just see whether it works for the rod. For an axis through one end, the moment of inertia should be $ML^2/3$, for we calculated that. The center of mass of a rod, of course, is in the center of the rod, at a distance $L/2$. Therefore we should find that $ML^2/3 = ML^2/12 + M(L/2)^2$. Since one-quarter plus one-twelfth is one-third, we have made no fundamental error.

Incidentally, we did not really need to use an integral to find the moment of inertia (19.5). If we simply assume that it is ML^2 times γ , an unknown coefficient, and then use the argument about the two halves to get $\frac{1}{4}\gamma$ for (19.6), then from our argument about transferring the axes we could prove that $\gamma = \frac{1}{4}\gamma + \frac{1}{4}$, so γ must be $1/3$. There is always another way to do it!

In applying the parallel-axis theorem, it is of course important to remember that the axis for I_c *must be parallel* to the axis about which the moment of inertia is wanted.

One further property of the moment of inertia is worth mentioning because it is often helpful in finding the moment of inertia of certain kinds of objects. This property is that if one has a *plane figure* and a set of coordinate axes with origin in the plane and z -axis perpendicular to the plane, then the moment of inertia of this figure about the z -axis is equal to the sum of the moments of inertia about the x - and y -axes. This is easily proved by noting that

$$I_z = \sum m_i(y_i^2 + z_i^2) = \sum m_i y_i^2$$

(since $z_i = 0$). Similarly,

$$I_y = \sum m_i(x_i^2 + z_i^2) = \sum m_i x_i^2,$$

but

$$\begin{aligned} I_z &= \sum m_i(x_i^2 + y_i^2) = \sum m_i x_i^2 + \sum m_i y_i^2 \\ &= I_x + I_y. \end{aligned}$$

As an example, the moment of inertia of a uniform rectangular plate of mass M , width w , and length L , about an axis perpendicular to the plate and through its center is simply

$$I = M(w^2 + L^2)/12,$$

because its moment of inertia about an axis in its plane and parallel to its length is $Mw^2/12$, i.e., just as for a rod of length w , and the moment of inertia about the other axis in its plane is $ML^2/12$, just as for a rod of length L .

To summarize, the moment of inertia of an object about a given axis, which we shall call the z -axis, has the following properties:

- (1) The moment of inertia is

$$I_z = \sum_i m_i(x_i^2 + y_i^2) = \int (x^2 + y^2) dm.$$

- (2) If the object is made of a number of parts, each of whose moment of inertia is known, the total moment of inertia is the sum of the moments of inertia of the pieces.
- (3) The moment of inertia about any given axis is equal to the moment of inertia about a parallel axis through the CM plus the total mass times the square of the distance from the axis to the CM.
- (4) If the object is a plane figure, the moment of inertia about an axis perpendicular to the plane is equal to the sum of the moments of inertia about any two mutually perpendicular axes lying in the plane and intersecting at the perpendicular axis.

The moments of inertia of a number of elementary shapes having uniform mass densities are given in Table 19-1, and the moments of inertia of some other objects, which may be deduced from Table 19-1, using the above properties, are given in Table 19-2.

Table 19-1

Object	z -axis	I_z
Thin rod, length L	\perp rod at center	$ML^2/12$
Thin concentric circular ring, radii r_1 and r_2	\perp ring at center	$M(r_1^2 + r_2^2)/2$
Sphere, radius r	through center	$2Mr^2/5$

Table 19-2

Object	z -axis	I_z
Rect. sheet, sides a, b	$\parallel b$ at center	$Ma^2/12$
Rect. sheet, sides a, b	\perp sheet at center	$M(a^2 + b^2)/12$
Thin annular ring, radii r_1, r_2	any diameter	$M(r_1^2 + r_2^2)/4$
Rect. parallelepiped, sides a, b, c	$\parallel c$, through center	$M(a^2 + b^2)/12$
Rt. circ. cyl., radius r , length L	$\parallel L$, through center	$Mr^2/2$
Rt. circ. cyl., radius r , length L	$\perp L$, through center	$M(r^2/4 + L^2/12)$

19-4 Rotational kinetic energy

Now let us go on to discuss dynamics further. In the analogy between linear motion and angular motion that we discussed in Chapter 18, we used the work theorem, but we did not talk about kinetic energy. What is the kinetic energy of a rigid body, rotating about a certain axis with an angular velocity ω ? We can immediately guess the correct answer by using our analogies. The moment of inertia corresponds to the mass, angular velocity corresponds to velocity, and so the kinetic energy ought to be $\frac{1}{2}I\omega^2$, and indeed it is, as will now be demonstrated. Suppose the object is rotating about some axis so that each point has a velocity whose magnitude is ωr_i , where r_i is the radius from the particular point to the axis.

Then if m_i is the mass of that point, the total kinetic energy of the whole thing is just the sum of the kinetic energies of all of the little pieces:

$$T = \frac{1}{2} \sum m_i v_i^2 = \frac{1}{2} \sum m_i (r_i \omega)^2.$$

Now ω^2 is a constant, the same for all points. Thus

$$T = \frac{1}{2} \omega^2 \sum m_i r_i^2 = \frac{1}{2} I \omega^2. \quad (19.8)$$

At the end of Chapter 18 we pointed out that there are some interesting phenomena associated with an object which is not rigid, but which changes from one rigid condition with a definite moment of inertia, to another rigid condition. Namely, in our example of the turntable, we had a certain moment of inertia I_1 with our arms stretched out, and a certain angular velocity ω_1 . When we pulled our arms in, we had a different moment of inertia, I_2 , and a different angular velocity, ω_2 , but again we were "rigid." The angular momentum remained constant, since there was no torque about the vertical axis of the turntable. This means that $I_1 \omega_1 = I_2 \omega_2$. Now what about the energy? That is an interesting question. With our arms pulled in, we turn faster, but our moment of inertia is less, and it looks as though the energies might be equal. But they are not, because what *does* balance is $I\omega$, not $I\omega^2$. So if we compare the kinetic energy before and after, the kinetic energy before is $\frac{1}{2} I_1 \omega_1^2 = \frac{1}{2} L \omega_1$, where $L = I_1 \omega_1 = I_2 \omega_2$ is the angular momentum. Afterward, by the same argument, we have $T = \frac{1}{2} L \omega_2$, and since $\omega_2 > \omega_1$ the kinetic energy of rotation is greater than it was before. So we had a certain energy when our arms were out, and when we pulled them in, we were turning faster and had more kinetic energy. What happened to the theorem of the conservation of energy? Somebody must have done some work. We did work! When did we do any work? When we move a weight horizontally, we do not do any work. If we hold a thing out and pull it in, we do not do any work. But that is when we are not rotating! When we *are* rotating, there is centrifugal force on the weights. They are trying to fly out, so when we are going around we have to pull the weights in against the centrifugal force. So, the work we do against the centrifugal force ought to agree with the difference in rotational energy, and of course it does. That is where the extra kinetic energy comes from.

There is still another interesting feature which we can treat only descriptively, as a matter of general interest. This feature is a little more advanced, but is worth pointing out because it is quite curious and produces many interesting effects.

Consider that turntable experiment again. Consider the body and the arms separately, from the point of view of the man who is rotating. After the weights are pulled in, the whole object is spinning faster, but observe, *the central part of the body is not changed*, yet it is turning faster after the event than before. So, if we were to draw a circle around the inner body, and consider only objects inside the circle, *their* angular momentum would *change*; they are going faster. Therefore there must be a torque exerted on the body while we pull in our arms. No torque can be exerted by the centrifugal force, because that is radial. So that means that among the forces that are developed in a rotating system, centrifugal force is not the entire story, *there is another force*. This other force is called *Coriolis force*, and it has the very strange property that when we move something in a rotating system, it seems to be pushed sidewise. Like the centrifugal force, it is an apparent force. But if we live in a system that is rotating, and move something radially, we find that we must also push it sidewise to move it radially. This sidewise push which we have to exert is what turned our body around.

Now let us develop a formula to show how this Coriolis force really works. Suppose Moe is sitting on a carousel that appears to him to be stationary. But from the point of view of Joe, who is standing on the ground and who knows the right laws of mechanics, the carousel is going around. Suppose that we have drawn a radial line on the carousel, and that Moe is moving some mass radially along this line. We would like to demonstrate that a sidewise force is required to do that. We can do this by paying attention to the angular momentum of the mass. It is

always going around with the same angular velocity ω , so that the angular momentum is

$$L = mv_{\text{tang}}r = m\omega r \cdot r = m\omega r^2.$$

So when the mass is close to the center, it has relatively little angular momentum, but if we move it to a new position farther out, if we increase r , m has more angular momentum, so *a torque must be exerted* in order to move it along the radius. (To walk along the radius in a carousel, one has to lean over and push sidewise. Try it sometime.) The torque that is required is the rate of change of L with time as m moves along the radius. If m moves only along the radius, ω stays constant, so that the torque is

$$\tau = F_c r = \frac{dL}{dt} = \frac{d(m\omega r^2)}{dt} = 2m\omega r \frac{dr}{dt},$$

where F_c is the Coriolis force. What we really want to know is what sidewise *force* has to be exerted by Moe in order to move m out at speed $v_r = dr/dt$. This is $F_c = \tau/r = 2m\omega v_r$.

Now that we have a formula for the Coriolis force, let us look at the situation a little more carefully, to see whether we can understand the origin of this force from a more elementary point of view. We note that the Coriolis force is the same at every radius, and is evidently present even at the origin! But it is especially easy to understand it at the origin, just by looking at what happens from the inertial system of Joe, who is standing on the ground. Figure 19-4 shows three successive views of m just as it passes the origin at $t = 0$. Because of the rotation of the carousel, we see that m does not move in a straight line, but in a *curved path* tangent to a diameter of the carousel where $r = 0$. In order for m to go in a curve, there must be a force to accelerate it in absolute space. This is the Coriolis force.

This is not the only case in which the Coriolis force occurs. We can also show that if an object is moving with constant speed around the circumference of a circle, there is also a Coriolis force. Why? Moe sees a velocity v_M around the circle. On the other hand, Joe sees m going around the circle with the velocity $v_J = v_M + \omega r$, because m is also carried by the carousel. Therefore we know what the force really is, namely, the total centripetal force due to the velocity v_J , or mv_J^2/r ; that is the actual force. Now from Moe's point of view, this centripetal force has three pieces. We may write it all out as follows:

$$F_r = -\frac{mv_J^2}{r} = -\frac{mv_M^2}{r} - 2mv_M\omega - m\omega^2 r.$$

Now, F_r is the force that Moe would see. Let us try to understand it. Would Moe appreciate the first term? "Yes," he would say, "even if I were not turning, there would be a centripetal force if I were to run around a circle with velocity v_M ." This is simply the centripetal force that Moe would expect, having nothing to do with rotation. In addition, Moe is quite aware that there is another centripetal force that would act even on objects which are standing still on his carousel. This is the third term. But there is another term in addition to these, namely the second term, which is again $2m\omega v$. The Coriolis force F_c was tangential when the velocity was radial, and now it is radial when the velocity is tangential. In fact, one expression has a minus sign relative to the other. The force is always in the same direction, relative to the velocity, no matter in which direction the velocity is. The force is at right angles to the velocity, and of magnitude $2m\omega v$.

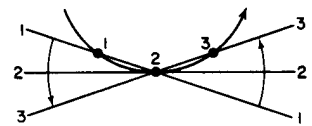


Fig. 19-4. Three successive views of a point moving radially on a rotating turntable.

Rotation in space

20-1 Torques in three dimensions

In this chapter we shall discuss one of the most remarkable and amusing consequences of mechanics, the behavior of a rotating wheel. In order to do this we must first extend the mathematical formulation of rotational motion, the principles of angular momentum, torque, and so on, to three-dimensional space. We shall not *use* these equations in all their generality and study all their consequences, because this would take many years, and we must soon turn to other subjects. In an introductory course we can present only the fundamental laws and apply them to a very few situations of special interest.

First, we notice that if we have a rotation in three dimensions, whether of a rigid body or any other system, what we deduced for two dimensions is still right. That is, it is still true that $xF_y - yF_x$ is the torque “in the xy -plane,” or the torque “around the z -axis.” It also turns out that this torque is still equal to the rate of change of $xp_y - yp_x$, for if we go back over the derivation of Eq. (18.15) from Newton’s laws we see that we did not have to assume that the motion was in a plane; when we differentiate $xp_y - yp_x$, we get $xF_y - yF_x$, so this theorem is still right. The quantity $xp_y - yp_x$, then, we call the angular momentum belonging to the xy -plane, or the angular momentum about the z -axis. This being true, we can use any other pair of axes and get another equation. For instance, we can use the yz -plane, and it is clear from symmetry that if we just substitute y for x and z for y , we would find $yF_z - zF_y$ for the torque and $yp_z - zp_y$ would be the angular momentum associated with the yz -plane. Of course we could have another plane, the zx -plane, and for this we would find $zF_x - xF_z = d/dt (zp_x - xp_z)$.

That these three equations can be deduced for the motion of a single particle is quite clear. Furthermore, if we added such things as $xp_y - yp_x$ together for many particles and called it the total angular momentum, we would have three kinds for the three planes xy , yz , and zx , and if we did the same with the forces, we would talk about the torque in the planes xy , yz , and zx also. Thus we would have laws that the external torque associated with any plane is equal to the rate of change of the angular momentum associated with that plane. This is just a generalization of what we wrote in two dimensions.

But now one may say, “Ah, but there are more planes; after all, can we not take some other plane at some angle, and calculate the torque on that plane from the forces? Since we would have to write another set of equations for every such plane, we would have a lot of equations!” Interestingly enough, it turns out that if we were to work out the combination $x'F_{y'} - y'F_{x'}$ for another plane, measuring the x' , $F_{y'}$, etc., in that plane, the result can be written as some *combination* of the three expressions for the xy -, yz - and zx -planes. There is nothing new. In other words, if we know what the three torques in the xy -, yz -, and zx -planes are, then the torque in any other plane, and correspondingly the angular momentum also, can be written as some combination of these: six percent of one and ninety-two percent of another, and so on. This property we shall now analyze.

Suppose that in the xyz -axes, Joe has worked out all his torques and his angular momenta in his planes. But Moe has axes x' , y' , z' in some other direction. To make it a little easier, we shall suppose that only the x - and y -axes have been turned. Moe’s x' and y' are new, but his z' happens to be the same. That is, he has new planes, let us say, for yz and zx . He therefore has new torques and angular momenta which he would work out. For example, his torque in the $x'y'$ -plane would be equal to $x'F_{y'} - y'F_{x'}$ and so forth. What we must now do is to find the relationship between the new torques and the old torques, so we will be able to make a

20-1 Torques in three dimensions

20-2 The rotation equations using cross products

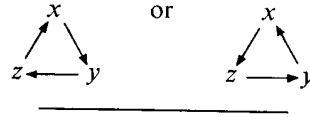
20-3 The gyroscope

20-4 Angular momentum of a solid body

connection from one set of axes to the other. Someone may say, "That looks just like what we did with vectors." And indeed, that is exactly what we are intending to do. Then he may say, "Well, isn't torque just a vector?" It *does* turn out to be a vector, but we do not know that right away without making an analysis. So in the following steps we shall make the analysis. We shall not discuss every step in detail, since we only want to illustrate how it works. The torques calculated by Joe are

$$\begin{aligned}\tau_{xy} &= xF_y - yF_x, \\ \tau_{yz} &= yF_z - zF_y, \\ \tau_{zx} &= zF_x - xF_z.\end{aligned}\tag{20.1}$$

We digress at this point to note that in such cases as this one may get the wrong sign for some quantity if the coordinates are not handled in the right way. Why not write $\tau_{yz} = zF_y - yF_z$? The problem arises from the fact that a coordinate system may be either "right-handed" or "left-handed." Having chosen (arbitrarily) a sign for, say τ_{xy} , then the correct expressions for the other two quantities may always be found by interchanging the letters xyz in either order



Moe now calculates the torques in his system:

$$\begin{aligned}\tau_{x'y'} &= x'F_{y'} - y'F_{x'}, \\ \tau_{y'z'} &= y'F_{z'} - z'F_{y'}, \\ \tau_{z'x'} &= z'F_{x'} - x'F_{z'}.\end{aligned}\tag{20.2}$$

Now we suppose that one coordinate system is rotated by a fixed angle θ , such that the z - and z' -axes are the same. (This angle θ has nothing to do with rotating objects or what is going on inside the coordinate system. It is merely the relationship between the axes used by one man and the axes used by the other, and is supposedly constant.) Thus the coordinates of the two systems are related by

$$\begin{aligned}x' &= x \cos \theta + y \sin \theta, \\ y' &= y \cos \theta - x \sin \theta, \\ z' &= z.\end{aligned}\tag{20.3}$$

Likewise, because force is a vector it transforms into the new system in the same way as do x , y , and z , since a thing is a vector if and only if the various components transform in the same way as x , y , and z :

$$\begin{aligned}F_{x'} &= F_x \cos \theta + F_y \sin \theta, \\ F_{y'} &= F_y \cos \theta - F_x \sin \theta, \\ F_{z'} &= F_z.\end{aligned}\tag{20.4}$$

Now we can find out how the torque transforms by merely substituting for x' , y' , and z' the expressions (20.3), and for $F_{x'}$, $F_{y'}$, $F_{z'}$ those given by (20.4), all into (20.2). So, we have a rather long string of terms for $\tau_{x'y'}$ and (rather surprisingly at first) it turns out that it comes right down to $xF_y - yF_x$, which we recognize to be the torque in the xy -plane:

$$\begin{aligned}\tau_{x'y'} &= (x \cos \theta + y \sin \theta)(F_y \cos \theta - F_x \sin \theta) \\ &\quad - (y \cos \theta - x \sin \theta)(F_x \cos \theta + F_y \sin \theta) \\ &= xF_y(\cos^2 \theta + \sin^2 \theta) - yF_x(\sin^2 \theta + \cos^2 \theta) \\ &\quad + xF_x(-\sin \theta \cos \theta + \sin \theta \cos \theta) \\ &\quad + yF_y(\sin \theta \cos \theta - \sin \theta \cos \theta) \\ &= xF_y - yF_x = \tau_{xy}.\end{aligned}\tag{20.5}$$

That result is clear, for if we only turn our axes *in the plane*, the twist around z in that plane is no different than it was before, because it is the same plane! What will be more interesting is the expression for $\tau_{y'z'}$, because that is a new plane. We now do exactly the same thing with the $y'z'$ -plane, and it comes out as follows:

$$\begin{aligned}\tau_{y'z'} &= (y \cos \theta - x \sin \theta)F_z \\ &\quad - z(F_y \cos \theta - F_x \sin \theta) \\ &= (yF_z - zF_y) \cos \theta + (zF_x - xF_z) \sin \theta \\ &= \tau_{yz} \cos \theta + \tau_{zx} \sin \theta.\end{aligned}\tag{20.6}$$

Finally, we do it for $z'x'$:

$$\begin{aligned}\tau_{z'x'} &= z(F_x \cos \theta + F_y \sin \theta) \\ &\quad - (x \cos \theta + y \sin \theta)F_z \\ &= (zF_x - xF_z) \cos \theta - (yF_z - zF_y) \sin \theta \\ &= \tau_{zx} \cos \theta - \tau_{yz} \sin \theta.\end{aligned}\tag{20.7}$$

We wanted to get a rule for finding torques in new axes in terms of torques in old axes, and now we have the rule. How can we ever remember that rule? If we look carefully at (20.5), (20.6), and (20.7), we see that there is a close relationship between these equations and the equations for x , y , and z . If, somehow, we could call τ_{xy} the z -component of something, let us call it the z -component of τ , then it would be all right; we would understand (20.5) as a vector transformation, since the z -component would be unchanged, as it should be. Likewise, if we associate with the yz -plane the x -component of our newly invented vector, and with the zx -plane, the y -component, then these transformation expressions would read

$$\begin{aligned}\tau_{z'} &= \tau_z, \\ \tau_{x'} &= \tau_x \cos \theta + \tau_y \sin \theta, \\ \tau_{y'} &= \tau_y \cos \theta - \tau_x \sin \theta,\end{aligned}\tag{20.8}$$

which is just the rule for vectors!

Therefore we have proved that we may identify the combination of $xF_y - yF_x$ with what we ordinarily call the z -component of a certain artificially invented vector. Although a torque is a twist on a plane, and it has no *a priori* vector character, mathematically it does behave like a vector. This vector is at right angles to the plane of the twist, and its length is proportional to the strength of the twist. The three components of such a quantity will transform like a real vector.

So we represent torques by vectors; with each plane on which the torque is supposed to be acting, we associate a line at right angles, by a rule. But "at right angles" leaves the sign unspecified. To get the sign right, we must adopt a rule which will tell us that if the torque were in a certain sense on the xy -plane, then the axis that we want to associate with it is in the "up" z -direction. That is, somebody has to define "right" and "left" for us. Supposing that the coordinate system is x , y , z in a right-hand system, then the rule will be the following: if we think of the twist as if we were turning a screw having a right-hand thread, then the direction of the vector that we will associate with that twist is in the direction that the screw would advance.

Why is torque a vector? It is a miracle of good luck that we can associate a single axis with a plane, and therefore that we can associate a vector with the torque; it is a special property of three-dimensional space. In two dimensions, the torque is an ordinary scalar, and there need be no direction associated with it. In three dimensions, it is a vector. If we had four dimensions, we would be in great difficulty, because (if we had time, for example, as the fourth dimension) we would not only have planes like xy , yz , and zx , we would also have tx -, ty -, and tz -planes. There would be *six* of them, and one cannot represent six quantities as one vector in four dimensions.

We will be living in three dimensions for a long time, so it is well to notice that the foregoing mathematical treatment did not depend upon the fact that x

was position and F was force; it only depended on the transformation laws for vectors. Therefore if, instead of x , we used the x -component of some other vector, it is not going to make any difference. In other words, if we were to calculate $a_x b_y - a_y b_x$, where \mathbf{a} and \mathbf{b} are vectors, and call it the z -component of some new quantity c , then these new quantities form a vector \mathbf{c} . We need a mathematical notation for the relationship of the new vector, with its three components, to the vectors \mathbf{a} and \mathbf{b} . The notation that has been devised for this is $\mathbf{c} = \mathbf{a} \times \mathbf{b}$. We have then, in addition to the ordinary scalar product in the theory of vector analysis, a new kind of product, called the *vector product*. Thus, if $\mathbf{c} = \mathbf{a} \times \mathbf{b}$, this is the same as writing

$$\begin{aligned} c_x &= a_y b_z - a_z b_y, \\ c_y &= a_z b_x - a_x b_z, \\ c_z &= a_x b_y - a_y b_x. \end{aligned} \quad (20.9)$$

If we reverse the order of \mathbf{a} and \mathbf{b} , calling \mathbf{a} , \mathbf{b} and \mathbf{b} , \mathbf{a} , we would have the sign of \mathbf{c} reversed, because c_z would be $b_x a_y - b_y a_x$. Therefore the cross product is unlike ordinary multiplication, where $ab = ba$; for the cross product, $\mathbf{b} \times \mathbf{a} = -\mathbf{a} \times \mathbf{b}$. From this, we can prove at once that if $\mathbf{a} = \mathbf{b}$, the cross product is zero. Thus, $\mathbf{a} \times \mathbf{a} = 0$.

The cross product is very important for representing the features of rotation, and it is important that we understand the geometrical relationship of the three vectors \mathbf{a} , \mathbf{b} , and \mathbf{c} . Of course the relationship in components is given in Eq. (20.9) and from that one can determine what the relationship is in geometry. The answer is, first, that the vector \mathbf{c} is perpendicular to both \mathbf{a} and \mathbf{b} . (Try to calculate $\mathbf{c} \cdot \mathbf{a}$, and see if it does not reduce to zero.) Second, the magnitude of \mathbf{c} turns out to be the magnitude of \mathbf{a} times the magnitude of \mathbf{b} times the sine of the angle between the two. In which direction does \mathbf{c} point? Imagine that we turn \mathbf{a} into \mathbf{b} through an angle less than 180° ; a screw with a right-hand thread turning in this way will advance in the direction of \mathbf{c} . The fact that we say a *right-hand* screw instead of a *left-hand* screw is a convention, and is a perpetual reminder that if \mathbf{a} and \mathbf{b} are “honest” vectors in the ordinary sense, the new kind of “vector” which we have created by $\mathbf{a} \times \mathbf{b}$ is artificial, or slightly different in its character from \mathbf{a} and \mathbf{b} , because it was made up with a special rule. If \mathbf{a} and \mathbf{b} are called ordinary vectors, we have a special name for them, we call them *polar vectors*. Examples of such vectors are the coordinate \mathbf{r} , force \mathbf{F} , momentum \mathbf{p} , velocity \mathbf{v} , electric field \mathbf{E} , etc.; these are ordinary polar vectors. Vectors which involve just one cross product in their definition are called *axial vectors* or *pseudovectors*. Examples of pseudovectors are, of course, torque $\boldsymbol{\tau}$ and the angular momentum \mathbf{L} . It also turns out that the angular velocity $\boldsymbol{\omega}$ is a pseudovector, as is the magnetic field \mathbf{B} .

In order to complete the mathematical properties of vectors, we should know all the rules for their multiplication, using dot and cross products. In our applications at the moment, we will need very little of this, but for the sake of completeness we shall write down all of the rules for vector multiplication so that we can use the results later. These are

$$\begin{aligned} (a) \quad & \mathbf{a} \times (\mathbf{b} + \mathbf{c}) = \mathbf{a} \times \mathbf{b} + \mathbf{a} \times \mathbf{c}, \\ (b) \quad & (\alpha \mathbf{a}) \times \mathbf{b} = \alpha(\mathbf{a} \times \mathbf{b}), \\ (c) \quad & \mathbf{a} \cdot (\mathbf{b} \times \mathbf{c}) = (\mathbf{a} \times \mathbf{b}) \cdot \mathbf{c}, \\ (d) \quad & \mathbf{a} \times (\mathbf{b} \times \mathbf{c}) = \mathbf{b}(\mathbf{a} \cdot \mathbf{c}) - \mathbf{c}(\mathbf{a} \cdot \mathbf{b}), \\ (e) \quad & \mathbf{a} \times \mathbf{a} = 0, \\ (f) \quad & \mathbf{a} \cdot (\mathbf{a} \times \mathbf{b}) = 0. \end{aligned} \quad (20.10)$$

20-2 The rotation equations using cross products

Now let us ask whether any equations in physics can be written using the cross product. The answer, of course, is that a great many equations can be so written. For instance, we see immediately that the torque is equal to the position

vector cross the force:

$$\boldsymbol{\tau} = \mathbf{r} \times \mathbf{F}. \quad (20.11)$$

This is a vector summary of the three equations $\tau_x = yF_z - zF_y$, etc. By the same token, the angular momentum vector, if there is only one particle present, is the distance from the origin multiplied by the vector momentum:

$$\mathbf{L} = \mathbf{r} \times \mathbf{p}. \quad (20.12)$$

For three-dimensional space rotation, the dynamical law analogous to the law $\mathbf{F} = d\mathbf{p}/dt$ of Newton, is that the torque vector is the rate of change with time of the angular momentum vector:

$$\boldsymbol{\tau} = d\mathbf{L}/dt. \quad (20.13)$$

If we sum (20.13) over many particles, the external torque on a system is the rate of change of the total angular momentum:

$$\boldsymbol{\tau}_{\text{ext}} = d\mathbf{L}_{\text{tot}}/dt. \quad (20.14)$$

Another theorem: If the total external torque is zero, then the total vector angular momentum of the system is a constant. This is called the law of *conservation of angular momentum*. If there is no torque on a given system, its angular momentum cannot change.

What about angular velocity? Is it a vector? We have already discussed turning a solid object about a fixed axis, but for a moment suppose that we are turning it simultaneously about *two* axes. It might be turning about an axis inside a box, while the box is turning about some other axis. The net result of such combined motions is that the object simply turns about some new axis! The wonderful thing about this new axis is that it can be figured out this way. If the rate of turning in the xy -plane is written as a vector in the z -direction whose length is equal to the rate of rotation in the plane, and if another vector is drawn in the y -direction, say, which is the rate of rotation in the zx -plane, then if we add these together as a vector, the magnitude of the result tells us how fast the object is turning, and the direction tells us in what plane, by the rule of the parallelogram. That is to say, simply, angular velocity *is* a vector, where we draw the magnitudes of the rotations in the three planes as projections at right angles to those planes.*

As a simple application of the use of the angular velocity vector, we may evaluate the power being expended by the torque acting on a rigid body. The power, of course, is the rate of change of work with time; in three dimensions, the power turns out to be $P = \boldsymbol{\tau} \cdot \boldsymbol{\omega}$.

All the formulas that we wrote for plane rotation can be generalized to three dimensions. For example, if a rigid body is turning about a certain axis with angular velocity $\boldsymbol{\omega}$, we might ask, "What is the velocity of a point at a certain radial position \mathbf{r} ?" We shall leave it as a problem for the student to show that the velocity of a particle in a rigid body is given by $\mathbf{v} = \boldsymbol{\omega} \times \mathbf{r}$, where $\boldsymbol{\omega}$ is the angular velocity and \mathbf{r} is the position. Also, as another example of cross products, we had a formula for Coriolis force, which can also be written using cross products: $\mathbf{F}_c = 2m\mathbf{v} \times \boldsymbol{\omega}$. That is, if a particle is moving with velocity \mathbf{v} in a coordinate system which is, in fact, rotating with angular velocity $\boldsymbol{\omega}$, and we want to think in terms of the rotating coordinate system, then we have to add the pseudoforce \mathbf{F}_c .

20-3 The gyroscope

Let us now return to the law of conservation of angular momentum. This law may be demonstrated with a rapidly spinning wheel, or gyroscope, as follows (see Fig. 20-1). If we sit on a swivel chair and hold the spinning wheel with its axis horizontal, the wheel has an angular momentum about the horizontal axis.

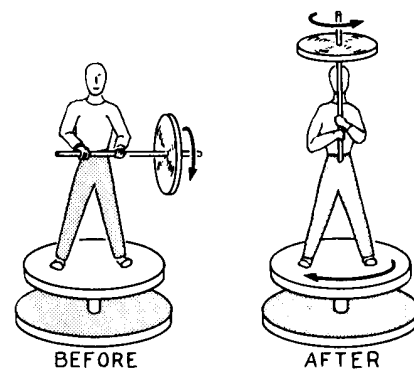


Fig. 20-1. *Before*: axis is horizontal; moment about vertical axis = 0. *After*: axis is vertical; momentum about vertical axis is still zero; man and chair spin in direction opposite to spin of the wheel.

* That this is true can be derived by compounding the displacements of the particles of the body during an infinitesimal time Δt . It is not self-evident, and is left to those who are interested to try to figure it out.

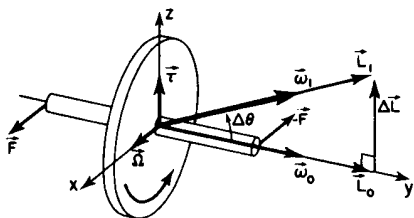


Fig. 20-2. A gyroscope.

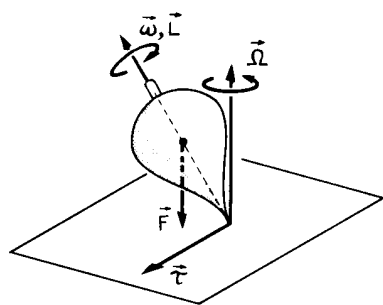


Fig. 20-3. A rapidly spinning top. Note that the direction of the torque vector is the direction of the precession.

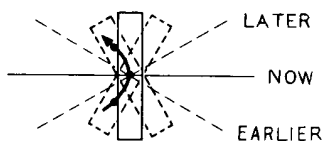


Fig. 20-4. The motion of particles in the spinning wheel of Fig. 20-2, whose axis is turning, is in curved lines.

Angular momentum around a *vertical* axis cannot change because of the (frictionless) pivot of the chair, so if we turn the axis of the wheel into the vertical, then the wheel would have angular momentum about the vertical axis, because it is now spinning about this axis. But the *system* (wheel, yourself, and chair) *cannot* have a vertical component, so we and the chair have to turn in the direction opposite to the spin of the wheel, to balance it.

First let us analyze in more detail the thing we have just described. What is surprising, and what we must understand, is the origin of the forces which turn us and the chair around as we turn the axis of the gyroscope toward the vertical. Figure 20-2 shows the wheel spinning rapidly about the *y*-axis. Therefore its angular velocity is about that axis and, it turns out, its angular momentum is likewise in that direction. Now suppose that we wish to rotate the wheel about the *x*-axis at a small angular velocity Ω ; what forces are required? After a short time Δt , the axis has turned to a new position, tilted at an angle $\Delta\theta$ with the horizontal. Since the major part of the angular momentum is due to the spin on the axis (very little is contributed by the slow turning), we see that the angular momentum vector has changed. What is the change in angular momentum? The angular momentum does not change in *magnitude*, but it does change in *direction* by an amount $\Delta\theta$. The magnitude of the vector ΔL is thus $\Delta L = L_0 \Delta\theta$, so that the torque, which is the time rate of change of the angular momentum, is $\tau = \Delta L / \Delta t = L_0 \Delta\theta / \Delta t = L_0 \Omega$. Taking the directions of the various quantities into account, we see that

$$\tau = \Omega \times L_0. \quad (20.15)$$

Thus, if Ω and L_0 are both horizontal, as shown in the figure, τ is *vertical*. To produce such a torque, horizontal forces F and $-F$ must be applied at the ends of the axle. How are these forces applied? By our hands, as we try to rotate the axis of the wheel into the vertical direction. But Newton's Third Law demands that equal and opposite forces (and equal and opposite *torques*) act on *us*. This causes us to rotate in the opposite sense about the vertical axis *z*.

This result can be generalized for a rapidly spinning top. In the familiar case of a spinning top, gravity acting on its center of mass furnishes a torque about the point of contact with the floor (see Fig. 20-3). This torque is in the horizontal direction, and causes the top to precess with its axis moving in a circular cone about the vertical. If Ω is the (vertical) angular velocity of precession, we again find that

$$\tau = d\mathbf{L}/dt = \Omega \times L_0.$$

Thus, when we apply a torque to a rapidly spinning top, the direction of the precessional motion is in the direction of the torque, or at right angles to the forces producing the torque.

We may now claim to understand the precession of gyroscopes, and indeed we do, mathematically. However, this is a mathematical thing which, in a sense, appears as a "miracle." It will turn out, as we go to more and more advanced physics, that many simple things can be deduced mathematically more rapidly than they can be really understood in a fundamental or simple sense. This is a strange characteristic, and as we get into more and more advanced work there are circumstances in which mathematics will produce results which *no one* has really been able to understand in any direct fashion. An example is the Dirac equation, which appears in a very simple and beautiful form, but whose consequences are hard to understand. In our particular case, the precession of a top looks like some kind of a miracle involving right angles and circles, and twists and right-hand screws. What we should try to do is to understand it in a more physical way.

How can we explain the torque in terms of the real forces and the accelerations? We note that when the wheel is precessing, the particles that are going around the wheel are not really moving in a plane because the wheel is precessing (see Fig. 20-4). As we explained previously (Fig. 19-4), the particles which are crossing through the precession axis are moving in *curved paths*, and this requires application of a lateral force. This is supplied by our pushing on the axle, which then com-

municates the force to the rim through the spokes. “Wait,” someone says, “what about the particles that are going back on the other side?” It does not take long to decide that there must be a force in the *opposite direction* on that side. The net force that we have to apply is therefore zero. The *forces* balance out, but one of them must be applied at one side of the wheel, and the other must be applied at the other side of the wheel. We could apply these forces directly, but because the wheel is solid we are allowed to do it by pushing on the axle, since forces can be carried up through the spokes.

What we have so far proved is that if the wheel is precessing, it can balance the torque due to gravity or some other applied torque. But all we have shown is that this is *a* solution of an equation. That is, if the torque is given, and *if we get the spinning started right*, then the wheel will precess smoothly and uniformly. But we have not proved (and it is not true) that a uniform precession is the *most general* motion a spinning body can undergo as the result of a given torque. The general motion involves also a “wobbling” about the mean precession. This “wobbling” is called *nutation*.

Some people like to say that when one exerts a torque on a gyroscope, it turns and it precesses, and that the torque *produces* the precession. It is very strange that when one suddenly lets go of a gyroscope, it does not *fall* under the action of gravity, but moves sidewise instead! Why is it that the *downward* force of the gravity, which we *know* and *feel*, makes it go *sidewise*? All the formulas in the world like (20.15) are not going to tell us, because (20.15) is a special equation, valid only after the gyroscope is precessing nicely. What really happens, in detail, is the following. If we were to hold the axis absolutely fixed, so that it cannot precess in any manner (but the top is spinning) then there is no torque acting, not even a torque from gravity, because it is balanced by our fingers. But if we suddenly let go, then there will instantaneously be a torque from gravity. Anyone in his right mind would think that the top would fall, and that is what it starts to do, as can be seen if the top is not spinning too fast.

The gyro actually does fall, as we would expect. But as soon as it falls, it is then turning, and if this turning were to continue, a torque would be required. In the absence of a torque in this direction, the gyro begins to “fall” in the direction opposite that of the missing force. This gives the gyro a component of motion around the vertical axis, as it would have in steady precession. But the actual motion “overshoots” the steady precessional velocity, and the axis actually rises again to the level from which it started. The path followed by the end of the axle is a cycloid (the path followed by a pebble that is stuck in the tread of an automobile tire). Ordinarily, this motion is too quick for the eye to follow, and it damps out quickly because of the friction in the gimbal bearings, leaving only the steady precessional drift (Fig. 20-5). The slower the wheel spins, the more obvious the nutation is.

When the motion settles down, the axis of the gyro is a little bit lower than it was at the start. Why? (These are the more complicated details, but we bring them in because we do not want the reader to get the idea that the gyroscope is an absolute miracle. It *is* a wonderful thing, but it is not a miracle.) If we were holding the axis absolutely horizontally, and suddenly let go, then the simple precession equation would tell us that it precesses, that it goes around in a horizontal plane. But that is impossible! Although we neglected it before, it is true that the wheel has *some* moment of inertia about the precession axis, and if it is moving about that axis, even slowly, it has a weak angular momentum about the axis. Where did it come from? If the pivots are perfect, there is no torque about the vertical axis. How then *does* it get to precess if there is no change in the angular momentum? The answer is that the cycloidal motion of the end of the axis damps down to the average, steady motion of the center of the equivalent rolling circle. That is, it settles down a little bit low. Because it is low, the spin angular momentum now has a small vertical component, which is exactly what is needed for the precession. So you see it has to go down a little, in order to go around. It has to yield a little bit to the gravity; by turning its axis down a little bit, it maintains the rotation about the vertical axis. That, then, is the way a gyroscope works.

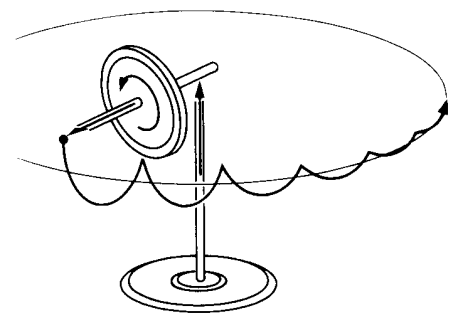


Fig. 20-5. Actual motion of tip of axis of gyroscope under gravity just after releasing axis previously held fixed.

20-4 Angular momentum of a solid body

Before we leave the subject of rotations in three dimensions, we shall discuss, at least qualitatively, a few effects that occur in three-dimensional rotations that are not self-evident. The main effect is that, in general, the angular momentum of a rigid body is *not necessarily* in the same direction as the angular velocity. Consider a wheel that is fastened onto a shaft in a lopsided fashion, but with the axis through the center of gravity, to be sure (Fig. 20-6). When we spin the wheel around the axis, anybody knows that there will be shaking at the bearings because of the lopsided way we have it mounted. Qualitatively, we know that in the rotating system there is centrifugal force acting on the wheel, trying to throw its mass as far as possible from the axis. This tends to line up the plane of the wheel so that it is perpendicular to the axis. To resist this tendency, a torque is exerted by the bearings. If there is a torque exerted by the bearings, there must be a rate of change of angular momentum. How can there be a rate of change of angular momentum when we are simply turning the wheel about the axis? Suppose we break the angular velocity ω into components ω_1 and ω_2 perpendicular and parallel to the plane of the wheel. What is the angular momentum? The moments of inertia about these two axes are *different*, so the angular momentum components, which (in these particular, special axes only) are equal to the moments of inertia times the corresponding angular velocity components, are in a *different ratio* than are the angular velocity components. Therefore the angular momentum vector is in a direction in space *not* along the axis. When we turn the object, we have to turn the angular momentum vector in space, so we must exert torques on the shaft.

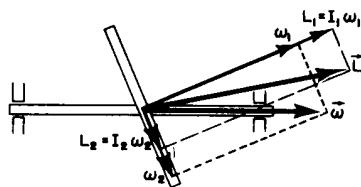


Fig. 20-6. The angular momentum of a rotating body is not necessarily parallel to the angular velocity.

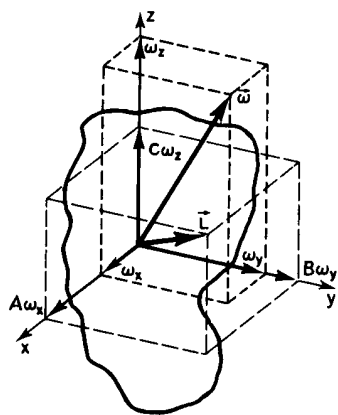


Fig. 20-7. The angular velocity and angular momentum of a rigid body ($A > B > C$).

Although it is much too complicated to prove here, there is a very important and interesting property of the moment of inertia which is easy to describe and to use, and which is the basis of our above analysis. This property is the following: Any rigid body, even an irregular one like a potato, possesses three mutually perpendicular axes through the CM, such that the moment of inertia about one of these axes has the greatest possible value for any axis through the CM, the moment of inertia about another of the axes has the *minimum* possible value, and the moment of inertia about the third axis is intermediate between these two (or equal to one of them). These axes are called the *principal axes* of the body, and they have the important property that if the body is rotating about one of them, its angular momentum is in the same direction as the angular velocity. For a body having axes of symmetry, the principal axes are along the symmetry axes.

If we take the x -, y -, and z -axes along the principal axes, and call the corresponding principal moments of inertia A , B , and C , we may easily evaluate the angular momentum and the kinetic energy of rotation of the body for any angular velocity ω . If we resolve ω into components ω_x , ω_y , and ω_z along the x -, y -, z -axes, and use unit vectors \mathbf{i} , \mathbf{j} , \mathbf{k} , also along x , y , z , we may write the angular momentum as

$$\mathbf{L} = A\omega_x\mathbf{i} + B\omega_y\mathbf{j} + C\omega_z\mathbf{k}. \quad (20.16)$$

The kinetic energy of rotation is

$$\begin{aligned} \text{KE} &= \frac{1}{2}(A\omega_x^2 + B\omega_y^2 + C\omega_z^2) \\ &= \frac{1}{2}\mathbf{L} \cdot \boldsymbol{\omega}. \end{aligned} \quad (20.17)$$

The Harmonic Oscillator

21-1 Linear differential equations

In the study of physics, usually the course is divided into a series of subjects, such as mechanics, electricity, optics, etc., and one studies one subject after the other. For example, this course has so far dealt mostly with mechanics. But a strange thing occurs again and again: the equations which appear in different fields of physics, and even in other sciences, are often almost exactly the same, so that many phenomena have analogs in these different fields. To take the simplest example, the propagation of sound waves is in many ways analogous to the propagation of light waves. If we study acoustics in great detail we discover that much of the work is the same as it would be if we were studying optics in great detail. So the study of a phenomenon in one field may permit an extension of our knowledge in another field. It is best to realize from the first that such extensions are possible, for otherwise one might not understand the reason for spending a great deal of time and energy on what appears to be only a small part of mechanics.

The harmonic oscillator, which we are about to study, has close analogs in many other fields; although we start with a mechanical example of a weight on a spring, or a pendulum with a small swing, or certain other mechanical devices, we are really studying a certain *differential equation*. This equation appears again and again in physics and in other sciences, and in fact it is a part of so many phenomena that its close study is well worth our while. Some of the phenomena involving this equation are the oscillations of a mass on a spring; the oscillations of charge flowing back and forth in an electrical circuit; the vibrations of a tuning fork which is generating sound waves; the analogous vibrations of the electrons in an atom, which generate light waves; the equations for the operation of a servosystem, such as a thermostat trying to adjust a temperature; complicated interactions in chemical reactions; the growth of a colony of bacteria in interaction with the food supply and the poisons the bacteria produce; foxes eating rabbits eating grass, and so on; all these phenomena follow equations which are very similar to one another, and this is the reason why we study the mechanical oscillator in such detail. The equations are called *linear differential equations with constant coefficients*. A linear differential equation with constant coefficients is a differential equation consisting of a sum of several terms, each term being a derivative of the dependent variable with respect to the independent variable, and multiplied by some constant. Thus

$$a_n d^n x/dt^n + a_{n-1} d^{n-1} x/dt^{n-1} + \cdots + a_1 dx/dt + a_0 x = f(t) \quad (21.1)$$

is called a linear differential equation of order n with constant coefficients (each a_i is constant).

21-2 The harmonic oscillator

Perhaps the simplest mechanical system whose motion follows a linear differential equation with constant coefficients is a mass on a spring: first the spring stretches to balance the gravity; once it is balanced, we then discuss the vertical displacement of the mass from its equilibrium position (Fig. 21-1). We shall call this upward displacement x , and we shall also suppose that the spring is perfectly linear, in which case the force pulling back when the spring is stretched is precisely proportional to the amount of stretch. That is, the force is $-kx$ (with a

21-1 Linear differential equations

21-2 The harmonic oscillator

21-3 Harmonic motion and circular motion

21-4 Initial conditions

21-5 Forced oscillations

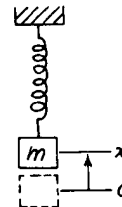


Fig. 21-1. A mass on a spring: a simple example of a harmonic oscillator.

minus sign to remind us that it pulls back). Thus the mass times the acceleration must equal $-kx$:

$$m d^2x/dt^2 = -kx. \quad (21.2)$$

For simplicity, suppose it happens (or we change our unit of time measurement) that the ratio $k/m = 1$. We shall first study the equation

$$d^2x/dt^2 = -x. \quad (21.3)$$

Later we shall come back to Eq. (21.2) with the k and m explicitly present.

We have already analyzed Eq. (21.3) in detail numerically; when we first introduced the subject of mechanics we solved this equation (see Eq. 9.12) to find the motion. By numerical integration we found a curve (Fig. 9-4) which showed that if m was initially displaced, but at rest, it would come down and go through zero; we did not then follow it any farther, but of course we know that it just keeps going up and down—it *oscillates*. When we calculated the motion numerically, we found that it went through the equilibrium point at $t = 1.570$. The length of the whole cycle is four times this long, or $t_0 = 6.28$ “sec.” This was found numerically, before we knew much calculus. We assume that in the meantime the Mathematics Department has brought forth a function which, when differentiated twice, is equal to itself with a minus sign. (There are, of course, ways of getting at this function in a direct fashion, but they are more complicated than already knowing what the answer is.) The function is $x = \cos t$. If we differentiate this we find $dx/dt = -\sin t$ and $d^2x/dt^2 = -\cos t = -x$. The function $x = \cos t$ starts, at $t = 0$, with $x = 1$, and no initial velocity; that was the situation with which we started when we did our numerical work. Now that we know that $x = \cos t$, we can calculate a *precise* value for the time at which it should pass $x = 0$. The answer is $t = \pi/2$, or 1.57108. We were wrong in the last figure because of the errors of numerical analysis, but it was very close!

Now to go further with the original problem, we restore the time units to real seconds. What is the solution then? First of all, we might think that we can get the constants k and m in by multiplying $\cos t$ by something. So let us try the equation $x = A \cos t$; then we find $dx/dt = -A \sin t$, and $d^2x/dt^2 = -A \cos t = -x$. Thus we discover to our horror that we did not succeed in solving Eq. (21.2), but we got Eq. (21.3) again! That fact illustrates one of the most important properties of linear differential equations: *if we multiply a solution of the equation by any constant, it is again a solution*. The mathematical reason for this is clear. If x is a solution, and we multiply both sides of the equation, say by A , we see that all derivatives are also multiplied by A , and therefore Ax is just as good a solution of the original equation as x was. The physics of it is the following. If we have a weight on a spring, and pull it down twice as far, the force is twice as much, the resulting acceleration is twice as great, the velocity it acquires in a given time is twice as great, the distance covered in a given time is twice as great; but it *has* to cover twice as great a distance in order to get back to the origin because it is pulled down twice as far. So it takes the *same time* to get back to the origin, irrespective of the initial displacement. In other words, with a linear equation, the motion has the same *time pattern*, no matter how “strong” it is.

That was the wrong thing to do—it only taught us that we can multiply the solution by anything, and it satisfies the same equation, but not a different equation. After a little cut and try to get to an equation with a different constant multiplying x , we find that we must alter the scale of *time*. In other words, Eq. (21.2) has a solution of the form

$$x = \cos \omega_0 t. \quad (21.4)$$

(It is important to realize that in the present case, ω_0 is not an angular velocity of a spinning body, but we run out of letters if we are not allowed to use the same letter for more than one thing.) The reason we put a subscript “0” on ω is that we are going to have more omegas before long; let us remember that ω_0 refers to the natural motion of this oscillator. Now we try Eq. (21.4) and this time we are more successful, because $dx/dt = -\omega_0 \sin \omega_0 t$ and $d^2x/dt^2 = -\omega_0^2 \cos \omega_0 t = -\omega_0^2 x$.

So at last we have solved the equation that we really wanted to solve. The equation $d^2x/dt^2 = -\omega_0^2x$ is the same as Eq. (21.2) if $\omega_0^2 = k/m$.

The next thing we must investigate is the physical significance of ω_0 . We know that the cosine function repeats itself when the angle it refers to is 2π . So $x = \cos \omega_0 t$ will repeat its motion, it will go through a complete cycle, when the "angle" changes by 2π . The quantity $\omega_0 t$ is often called the *phase* of the motion. In order to change $\omega_0 t$ by 2π , the time must change by an amount t_0 , called the *period* of one complete oscillation; of course t_0 must be such that $\omega_0 t_0 = 2\pi$. That is, $\omega_0 t_0$ must account for one cycle of the angle, and then everything will repeat itself—if we increase t by t_0 , we add 2π to the phase. Thus

$$t_0 = 2\pi/\omega_0 = 2\pi\sqrt{m/k}. \quad (21.5)$$

Thus if we had a heavier mass, it would take longer to oscillate back and forth on a spring. That is because it has more inertia, and so, while the forces are the same, it takes longer to get the mass moving. Or, if the spring is stronger, it will move more quickly, and that is right: the period is less if the spring is stronger.

Note that the period of oscillation of a mass on a spring does not depend in any way on *how* it has been started, how far down we pull it. The *period* is determined, but the amplitude of the oscillation is *not* determined by the equation of motion (21.2). The amplitude *is* determined, in fact, by how we let go of it, by what we call the *initial conditions* or starting conditions.

Actually, we have not quite found the most general possible solution of Eq. (21.2). There are other solutions. It should be clear why: because all of the cases covered by $x = a \cos \omega_0 t$ start with an initial displacement and no initial velocity. But it is possible, for instance, for the mass to start at $x = 0$, and we may then give it an impulsive kick, so that it has some speed at $t = 0$. Such a motion is not represented by a cosine—it is represented by a sine. To put it another way, if $x = \cos \omega_0 t$ is a solution, then is it not obvious that if we were to happen to walk into the room at some time (which we would call " $t = 0$ ") and saw the mass as it was passing $x = 0$, it would keep on going just the same? Therefore, $x = \cos \omega_0 t$ cannot be the most general solution; it must be possible to shift the beginning of time, so to speak. As an example, we could write the solution this way: $x = a \cos \omega_0(t - t_1)$, where t_1 is some constant. This also corresponds to shifting the origin of time to some new instant. Furthermore, we may expand

$$\cos(\omega_0 t + \Delta) = \cos \omega_0 t \cos \Delta - \sin \omega_0 t \sin \Delta,$$

and write

$$x = A \cos \omega_0 t + B \sin \omega_0 t,$$

where $A = a \cos \Delta$ and $B = -a \sin \Delta$. Any one of these forms is a possible way to write the complete, general solution of (21.2): that is, every solution of the differential equation $d^2x/dt^2 = -\omega_0^2x$ that exists in the world can be written as

$$(a) \quad x = a \cos \omega_0(t - t_1),$$

or

$$(b) \quad x = a \cos(\omega_0 t + \Delta), \quad (21.6)$$

or

$$(c) \quad x = A \cos \omega_0 t + B \sin \omega_0 t.$$

Some of the quantities in (21.6) have names: ω_0 is called the *angular frequency*; it is the number of radians by which the phase changes in a second. That is determined by the differential equation. The other constants are not determined by the equation, but by how the motion is started. Of these constants, a measures the maximum displacement attained by the mass, and is called the *amplitude* of oscillation. The constant Δ is sometimes called the *phase* of the oscillation, but that is a confusion, because other people call $\omega_0 t + \Delta$ the phase, and say the phase changes with time. We might say that Δ is a *phase shift* from some defined zero. Let us put it differently. Different Δ 's correspond to motions in different phases. That is true, but whether we want to call Δ *the* phase, or not, is another question.

21-3 Harmonic motion and circular motion

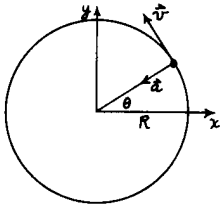


Fig. 21-2. A particle moving in a circular path at constant speed.

The fact that cosines are involved in the solution of Eq. (21.2) suggests that there might be some relationship to circles. This is artificial, of course, because there is no circle actually involved in the linear motion—it just goes up and down. We may point out that we have, in fact, already solved that differential equation when we were studying the mechanics of circular motion. If a particle moves in a circle with a constant speed v , the radius vector from the center of the circle to the particle turns through an angle whose size is proportional to the time. If we call this angle $\theta = vt/R$ (Fig. 21-2) then $d\theta/dt = \omega_0 = v/R$. We know that there is an acceleration $a = v^2/R = \omega_0^2 R$ toward the center. Now we also know that the position x , at a given moment, is the radius of the circle times $\cos \theta$, and that y is the radius times $\sin \theta$:

$$x = R \cos \theta, \quad y = R \sin \theta.$$

Now what about the acceleration? What is the x -component of acceleration, d^2x/dt^2 ? We have already worked that out geometrically; it is the magnitude of the acceleration times the cosine of the projection angle, with a minus sign because it is toward the center.

$$a_x = -a \cos \theta = -\omega^2 R \cos \theta = -\omega^2 x. \quad (21.7)$$

In other words, when a particle is moving in a circle, the horizontal component of its motion has an acceleration which is proportional to the horizontal displacement from the center. Of course we also have the solution for motion in a circle: $x = R \cos \omega_0 t$. Equation (21.7) does not depend upon the radius of the circle, so for a circle of any radius, one finds the same equation for a given ω_0 . Thus, for several reasons, we expect that the displacement of a mass on a spring will turn out to be proportional to $\cos \omega_0 t$, and will, in fact, be exactly the same motion as we would see if we looked at the x -component of the position of an object rotating in a circle with angular velocity ω_0 . As a check on this, one can devise an experiment to show that the up-and-down motion of a mass on a spring is the same as that of a point going around in a circle. In Fig. 21-3 an arc light projected on a screen casts shadows of a crank pin on a shaft and of a vertically oscillating mass, side by side. If we let go of the mass at the right time from the right place, and if the shaft speed is carefully adjusted so that the frequencies match, each should follow the other exactly. One can also check the numerical solution we obtained earlier with the cosine function, and see whether that agrees very well.

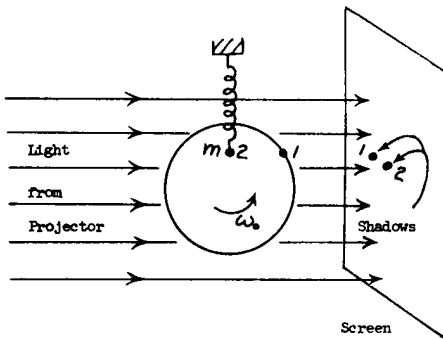


Fig. 21-3. Demonstration of the equivalence between simple harmonic motion and uniform circular motion.

Here we may point out that because uniform motion in a circle is so closely related mathematically to oscillatory up-and-down motion, we can analyze oscillatory motion in a simpler way if we imagine it to be a projection of something going in a circle. In other words, although the distance y means nothing in the oscillator problem, we may still artificially supplement Eq. (21.2) with another equation using y , and put the two together. If we do this, we will be able to analyze our one-dimensional oscillator with *circular* motions, which is a lot easier than having to solve a differential equation. The trick in doing this is to use complex numbers, a procedure we shall introduce in the next chapter.

21-4 Initial conditions

Now let us consider what determines the constants A and B , or a and Δ . Of course these are determined by how we start the motion. If we start the motion with just a small displacement, that is one type of oscillation; if we start with an initial displacement and then push up when we let go, we get still a different motion. The constants A and B , or a and Δ , or any other way of putting it, are determined, of course, by the way the motion started, not by any other features of the situation. These are called the *initial conditions*. We would like to connect the initial conditions with the constants. Although this can be done using any one of the forms (21.6), it turns out to be easiest if we use Eq. (21.6c). Suppose that at $t = 0$ we have started with an initial displacement x_0 and a certain velocity v_0 .

This is the most general way we can start the motion. (We cannot specify the *acceleration* with which it started, true, because that is determined by the spring, once we specify x_0 .) Now let us calculate A and B . We start with the equation for x ,

$$x = A \cos \omega_0 t + B \sin \omega_0 t.$$

Since we shall later need the velocity also, we differentiate x and obtain

$$v = -\omega_0 A \sin \omega_0 t + \omega_0 B \cos \omega_0 t.$$

These expressions are valid for all t , but we have special knowledge about x and v at $t = 0$. So if we put $t = 0$ into these equations, on the left we get x_0 and v_0 , because that is what x and v are at $t = 0$; also, we know that the cosine of zero is unity, and the sine of zero is zero. Therefore we get

$$x_0 = A \cdot 1 + B \cdot 0 = A$$

and

$$v_0 = -\omega_0 A \cdot 0 + \omega_0 B \cdot 1 = \omega_0 B.$$

So for this particular case we find that

$$A = x_0, \quad B = v_0/\omega_0.$$

From these values of A and B , we can get a and Δ if we wish.

That is the end of our solution, but there is one physically interesting thing to check, and that is the conservation of energy. Since there are no frictional losses, energy ought to be conserved. Let us use the formula

$$x = a \cos (\omega_0 t + \Delta);$$

then

$$v = -\omega_0 a \sin (\omega_0 t + \Delta).$$

Now let us find out what the kinetic energy T is, and what the potential energy U is. The potential energy at any moment is $\frac{1}{2}kx^2$, where x is the displacement and k is the constant of the spring. If we substitute for x , using our expression above, we get

$$U = \frac{1}{2}kx^2 = \frac{1}{2}ka^2 \cos^2 (\omega_0 t + \Delta).$$

Of course the potential energy is not constant; the potential never becomes negative, naturally—there is always some energy in the spring, but the amount of energy fluctuates with x . The kinetic energy, on the other hand, is $\frac{1}{2}mv^2$, and by substituting for v we get

$$T = \frac{1}{2}mv^2 = \frac{1}{2}m\omega_0^2 a^2 \sin^2 (\omega_0 t + \Delta).$$

Now the kinetic energy is zero when x is at the maximum, because then there is no velocity; on the other hand, it is maximal when x is passing through zero, because then it is moving fastest. This variation of the kinetic energy is just the opposite of that of the potential energy. But the total energy ought to be a constant. If we note that $k = m\omega_0^2$, we see that

$$T + U = \frac{1}{2}m\omega_0^2 a^2 [\cos^2 (\omega_0 t + \Delta) + \sin^2 (\omega_0 t + \Delta)] = \frac{1}{2}m\omega_0^2 a^2.$$

The energy is dependent on the square of the amplitude; if we have twice the amplitude, we get an oscillation which has four times the energy. The *average* potential energy is half the maximum and, therefore, half the total, and the average kinetic energy is likewise half the total energy.

21-5 Forced oscillations

Next we shall discuss the *forced harmonic oscillator*, i.e., one in which there is an external driving force acting. The equation then is the following:

$$m d^2 x / dt^2 = -kx + F(t). \quad (21.8)$$

We would like to find out what happens in these circumstances. The external driving force can have various kinds of functional dependence on the time; the first one that we shall analyze is very simple—we shall suppose that the force is oscillating:

$$F(t) = F_0 \cos \omega t. \quad (21.9)$$

Notice, however, that this ω is not necessarily ω_0 : we have ω under our control; the forcing may be done at different frequencies. So we try to solve Eq. (21.8) with the special force (21.9). What is the solution of (21.8)? One special solution, (we shall discuss the more general cases later) is

$$x = C \cos \omega t, \quad (21.10)$$

where the constant is to be determined. In other words, we might suppose that if we kept pushing back and forth, the mass would follow back and forth in step with the force. We can try it anyway. So we put (21.10) into (21.9), and get

$$-m\omega^2 C \cos \omega t = -m\omega_0^2 C \cos \omega t + F_0 \cos \omega t. \quad (21.11)$$

We have also put in $k = m\omega_0^2$, so that we will understand the equation better at the end. Now because the cosine appears everywhere, we can divide it out, and that shows that (21.10) is, in fact, a solution, provided we pick C just right. The answer is that C must be

$$C = F_0 / m(\omega_0^2 - \omega^2). \quad (21.12)$$

That is, m oscillates at the same frequency as the force, but with an amplitude which depends on the frequency of the force, and also upon the frequency of the natural motion of the oscillator. It means, first, that if ω is very small compared with ω_0 , then the displacement and the force are in the same direction. On the other hand, if we shake it back and forth very fast, then (21.12) tells us that C is negative if ω is above the natural frequency ω_0 of the harmonic oscillator. (We will call ω_0 the natural frequency of the harmonic oscillator, and ω the applied frequency.) At very high frequency the denominator may become very large, and there is then not much amplitude.

Of course the solution we have found is the solution only if things are started just right, for otherwise there is a part which usually dies out after a while. This other part is called the *transient* response to $F(t)$, while (21.10) and (21.12) are called the *steady-state* response.

According to our formula (21.12), a very remarkable thing should also occur: if ω is almost exactly the same as ω_0 , then C should approach infinity. So if we adjust the frequency of the force to be “in time” with the natural frequency, then we should get an enormous displacement. This is well known to anybody who has pushed a child on a swing. It does not work very well to close our eyes and push at a certain speed at random. If we happen to get the right timing, then the swing goes very high, but if we have the wrong timing, then sometimes we may be pushing when we should be pulling, and so on, and it does not work.

If we make ω exactly equal to ω_0 , we find that it should oscillate at an *infinite* amplitude, which is, of course, impossible. The reason it does not is that something goes wrong with the equation, there are some other frictional terms, and other forces, which are not in (21.8) but which occur in the real world. So the amplitude does not reach infinity for some reason; it may be that the spring breaks!

Algebra

22-1 Addition and multiplication

In our study of oscillating systems we shall have occasion to use one of the most remarkable, almost astounding, formulas in all of mathematics. From the physicist's point of view we could bring forth this formula in two minutes or so, and be done with it. But science is as much for intellectual enjoyment as for practical utility, so instead of just spending a few minutes on this amazing jewel, we shall surround the jewel by its proper setting in the grand design of that branch of mathematics which is called elementary algebra.

Now you may ask, "What is mathematics doing in a physics lecture?" We have several possible excuses: first, of course, mathematics is an important tool, but that would only excuse us for giving the formula in two minutes. On the other hand, in theoretical physics we discover that all our laws can be written in mathematical form; and that this has a certain simplicity and beauty about it. So, ultimately, in order to understand nature it may be necessary to have a deeper understanding of mathematical relationships. But the real reason is that the subject is enjoyable, and although we humans cut nature up in different ways, and we have different courses in different departments, such compartmentalization is really artificial, and we should take our intellectual pleasures where we find them.

Another reason for looking more carefully at algebra now, even though most of us studied algebra in high school, is that that was the first time we studied it; all the equations were unfamiliar, and it was hard work, just as physics is now. Every so often it is a great pleasure to look back to see what territory has been covered, and what the great map or plan of the whole thing is. Perhaps some day somebody in the Mathematics Department will present a lecture on mechanics in such a way as to show what it was we were trying to learn in the physics course!

The subject of algebra will not be developed from the point of view of a mathematician, exactly, because the mathematicians are mainly interested in how various mathematical facts are demonstrated, and how many assumptions are absolutely required, and what is not required. They are not so interested in the result of what they prove. For example, we may find the Pythagorean theorem quite interesting, that the sum of the squares of the sides of a right triangle is equal to the square of the hypotenuse; that is an interesting fact, a curiously simple thing, which may be appreciated without discussing the question of how to prove it, or what axioms are required. So, in the same spirit, we shall describe qualitatively, if we may put it that way, the system of elementary algebra. We say *elementary* algebra because there is a branch of mathematics called *modern* algebra in which some of the rules such as $ab = ba$, are abandoned, and it is still called algebra, but we shall not discuss that.

To discuss this subject we start in the middle. We suppose that we already know what integers are, what zero is, and what it means to increase a number by one unit. You may say, "That is not in the middle!" But it is the middle from a mathematical standpoint, because we could go even further back and describe the theory of sets in order to *derive* some of these properties of integers. But we are not going in that direction, the direction of mathematical philosophy and mathematical logic, but rather in the other direction, from the assumption that we know what integers are and we know how to count.

If we start with a certain number a , an integer, and we count successively one unit b times, the number we arrive at we call $a + b$, and that defines *addition* of integers.

22-1 Addition and multiplication

22-2 The inverse operations

22-3 Abstraction and generalization

22-4 Approximating irrational numbers

22-5 Complex numbers

22-6 Imaginary exponents

Once we have defined addition, then we can consider this: if we start with nothing and add a to it, b times in succession, we call the result *multiplication of integers*; we call it b times a .

Now we can also have a *succession of multiplications*: if we start with 1 and multiply by a , b times in succession, we call that *raising to a power*: a^b .

Now as a consequence of these definitions it can be easily shown that all of the following relationships are true:

$$\begin{array}{ll}
 \text{(a)} & a + b = b + a \\
 \text{(b)} & a + (b + c) = (a + b) + c \\
 \text{(c)} & ab = ba \\
 \text{(d)} & a(b + c) = ab + ac \\
 \text{(e)} & (ab)c = a(bc) \\
 \text{(f)} & (ab)^c = a^c b^c \\
 \text{(g)} & a^b a^c = a^{(b+c)} \\
 \text{(h)} & (a^b)^c = a^{(bc)} \\
 \text{(i)} & a + 0 = a \\
 \text{(j)} & a \cdot 1 = a \\
 \text{(k)} & a^1 = a
 \end{array} \tag{22.1}$$

These results are well known and we shall not belabor the point, we merely list them. Of course, 1 and 0 have special properties; for example, $a + 0$ is a , a times 1 = a , and a to the first power is a .

In this discussion we must also assume a few other properties like continuity and ordering, which are very hard to define; we will let the rigorous theory do it. Furthermore, it is definitely true that we have written down too many “rules”; some of them may be deducible from the others, but we shall not worry about such matters.

22-2 The inverse operations

In addition to the direct operations of addition, multiplication, and raising to a power, we have also the *inverse* operations, which are defined as follows. Let us assume that a and c are given, and that we wish to find what values of b satisfy such equations as $a + b = c$, $ab = c$, $b^a = c$. If $a + b = c$, b is defined as $c - a$, which is called *subtraction*. The operation called division is also clear: if $ab = c$, then $b = c/a$ defines division—a solution of the equation $ab = c$ “backwards.” Now if we have a power $b^a = c$ and we ask ourselves, “What is b ?” it is called the a th *root* of c : $b = \sqrt[a]{c}$. For instance, if we ask ourselves the following question, “What integer, raised to the third power, equals 8?” then the answer is called the *cube root* of 8; it is 2. Because b^a and a^b are not equal, there are *two* inverse problems associated with powers, and the other inverse problem would be, “To what power must we raise 2 to get 8?” This is called taking the *logarithm*. If $a^b = c$, we write $b = \log_a c$. The fact that it has a cumbersome notation relative to the others does not mean that it is any less elementary, at least applied to integers, than the other processes. Although logarithms come late in an algebra class, in practice they are, of course, just as simple as roots; they are just a different kind of solution of an algebraic equation. The direct and inverse operations are summarized as follows:

$$\begin{array}{ll}
 \text{(a)} & \text{addition} \\
 & a + b = c \\
 \text{(b)} & \text{multiplication} \\
 & ab = c \\
 \text{(c)} & \text{power} \\
 & b^a = c \\
 \text{(d)} & \text{power} \\
 & a^b = c \\
 \text{(a')} & \text{subtraction} \\
 & b = c - a \\
 \text{(b')} & \text{division} \\
 & b = c/a \\
 \text{(c')} & \text{root} \\
 & b = \sqrt[a]{c} \\
 \text{(d')} & \text{logarithm} \\
 & b = \log_a c
 \end{array} \tag{22.2}$$

Now here is the idea. These relationships, or rules, are correct for integers, since they follow from the definitions of addition, multiplication, and raising to a power. *We are going to discuss whether or not we can broaden the class of objects*

which a , b , and c represent so that they will obey these same rules, although the processes for $a + b$, and so on, will not be definable in terms of the direct action of adding 1, for instance, or successive multiplications by integers.

22-3 Abstraction and generalization

When we try to solve simple algebraic equations using all these definitions, we soon discover some insoluble problems, such as the following. Suppose that we try to solve the equation $b = 3 - 5$. That means, according to our definition of subtraction, that we must find a number which, when added to 5, gives 3. And of course there *is* no such number, because we consider only positive integers; this is an insoluble problem. However, the plan, the great idea, is this: *abstraction and generalization*. From the whole structure of algebra, rules plus integers, we abstract the original definitions of addition and multiplication, but we leave the rules (22.1) and (22.2), and assume these to be true *in general* on a wider class of numbers, even though they are originally derived on a smaller class. Thus, rather than using integers symbolically to define the rules, we use the rules as the definition of the symbols, which then represent a more general kind of number. As an example, by working with the rules alone we can show that $3 - 5 = 0 - 2$. In fact we can show that one can make *all* subtractions, provided we define a whole set of new numbers: $0 - 1$, $0 - 2$, $0 - 3$, $0 - 4$, and so on, called the *negative integers*. Then we may use all the other rules, like $a(b + c) = ab + ac$ and so forth, to find what the rules are for multiplying negative numbers, and we will discover, in fact, that all of the rules can be maintained with negative as well as positive integers.

So we have increased the range of objects over which the rules work, but the meaning of the symbols is different.

One cannot say, for instance, that -2 times 5 really means to add 5 together successively -2 times. That means nothing. But nevertheless everything will work out all right according to the rules.

An interesting problem comes up in taking powers. Suppose that we wish to discover what $a^{(3-5)}$ means. We know only that $3 - 5$ is a solution of the problem, $(3 - 5) + 5 = 3$. Knowing that, we know that $a^{(3-5)}a^5 = a^3$. Therefore $a^{(3-5)} = a^3/a^5$, by the definition of division. With a little more work, this can be reduced to $1/a^2$. So we find that the negative powers are the reciprocals of the positive powers, but $1/a^2$ is a meaningless symbol, because if a is a positive or negative integer, the square of it is greater than 1, and we do not yet know what we mean by 1 divided by a number greater than 1!

Onward! The great plan is to continue the process of generalization; whenever we find another problem that we cannot solve we extend our realm of numbers. Consider division: we cannot find a number which is an integer, even a negative integer, which is equal to the result of dividing 3 by 5. But if we suppose that all fractional numbers also satisfy the rules, then we can talk about multiplying and adding fractions, and everything works as well as it did before.

Take another example of powers: what is $a^{3/5}$? We know only that $(3/5)5 = 3$, since that was the definition of $3/5$. So we know also that $(a^{(3/5)})^5 = a^{(3/5)(5)} = a^3$, because this is one of the rules. Then by the definition of roots we find that $a^{(3/5)} = \sqrt[5]{a^3}$.

In this way, then, we can define what we mean by putting fractions in the various symbols, by using the rules themselves to help us determine the definition—it is not arbitrary. It is a remarkable fact that all the rules still work for positive and negative integers, as well as for fractions!

We go on in the process of generalization. Are there any other equations we cannot solve? Yes, there are. For example, it is impossible to solve this equation: $b = 2^{1/2} = \sqrt{2}$. It is impossible to find a number which is rational (a fraction) whose square is equal to 2. It is very easy for us in modern days to answer this question. We know the decimal system, and so we have no difficulty in appreciating the meaning of an unending decimal as a type of approximation to the square root of 2. Historically, this idea presented great difficulty to the Greeks. To really

define *precisely* what is meant here requires that we add some substance of continuity and ordering, and it is, in fact, quite the most difficult step in the processes of generalization just at this point. It was made, formally and rigorously, by Dedekind. However, without worrying about the mathematical rigor of the thing, it is quite easy to understand that what we mean is that we are going to find a whole sequence of approximate fractions, perfect fractions (because any decimal, when stopped somewhere, is of course rational), which just keeps on going, getting closer and closer to the desired result. That is good enough for what we wish to discuss, and it permits us to involve ourselves in irrational numbers, and to calculate things like the square root of 2 to any accuracy that we desire, with enough work.

22-4 Approximating irrational numbers

The next problem comes with what happens with the irrational powers. Suppose that we want to define, for instance, $10^{\sqrt{2}}$. In principle, the answer is simple enough. If we approximate the square root of 2 to a certain number of decimal places, then the power is rational, and we can take the approximate root, using the above method, and get an *approximation* to $10^{\sqrt{2}}$. Then we may run it up a few more decimal places (it is again rational), take the appropriate root, this time a much higher root because there is a much bigger denominator in the fraction, and get a better approximation. Of course we are going to get some enormously high roots involved here, and the work is quite difficult. How can we cope with this problem?

In the computations of square roots, cube roots, and other small roots, there is an arithmetical process available by which we can get one decimal place after another. But the amount of labor needed to calculate irrational powers and the logarithms that go with them (the inverse problem) is so great that there is no simple arithmetical process we can use. Therefore tables have been built up which permit us to calculate these powers, and these are called the tables of logarithms, or the tables of powers, depending on which way the table is set up. It is merely a question of saving time; if we must raise some number to an irrational power, we can look it up rather than having to compute it. Of course, such a computation is just a technical problem, but it is an interesting one, and of great historical value. In the first place, not only do we have the problem of solving $x = 10^{\sqrt{2}}$, but we also have the problem of solving $10^x = 2$, or $x = \log_{10} 2$. This is not a problem where we have to define a new kind of number for the result, it is merely a computational problem. The answer is simply an irrational number, an unending decimal, not a new kind of a number.

Let us now discuss the problem of calculating solutions of such equations. The general idea is really very simple. If we could calculate 10^1 , and $10^{4/10}$, and $10^{1/100}$, and $10^{4/1000}$ and so on, and multiply them all together, we would get $10^{1.414...}$ or $10^{\sqrt{2}}$, and that is the general idea on which things work. But instead of calculating $10^{1/10}$ and so on, we shall calculate $10^{1/2}$, $10^{1/4}$, and so on. Before we start, we should explain why we make so much work with 10, instead of some other number. Of course, we realize that logarithm tables are of great practical utility, quite aside from the mathematical problem of taking roots, since with any base at all,

$$\log_b (ac) = \log_b a + \log_b c. \quad (22.3)$$

We are all familiar with the fact that one can use this fact in a practical way to multiply numbers if we have a table of logarithms. The only question is, with what base b shall we compute? It makes no difference what base is used; we can use the same principle all the time, and if we are using logarithms to any particular base, we can find logarithms to any other base merely by a change in scale, a multiplying factor. If we multiply Eq. (22.3) by 61, it is just as true, and if we had a table of logs with a base b , and somebody else multiplied all of our table by 61, there would be no essential difference. Suppose that we know the logarithms of all the numbers to the base b . In other words, we can solve the equation $b^a = c$ for any c because

we have a table. The problem is to find the logarithm of the same number c to some other base, let us say the base x . We would like to solve $x^{a'} = c$. It is easy to do, because we can always write $x = b^t$, which defines t , knowing x and b . As a matter of fact, $t = \log_b x$. Then if we put that in and solve for a' , we see that $(b^t)^{a'} = b^{ta'} = c$. In other words, ta' is the logarithm of c in base b . Thus $a' = a/t$. Thus logs to base x are just $1/t$, which is a constant, times the logs to the base, b . Therefore any log table is equivalent to any other log table if we multiply by a constant, and the constant is $1/\log_b x$. This permits us to choose a particular base, and for convenience we take the base 10. (The question may arise as to whether there is any natural base, any base in which things are somehow simpler, and we shall try to find an answer to that later. At the moment we shall just use the base 10.)

Table 22-1
Successive Square Roots of Ten

Power s	$1024\ s$	10^s	$(10^s - 1)/s$
1	1024	10.00000	9.00
1/2	512	3.16228	4.32
1/4	256	1.77828	3.113
1/8	128	1.33352	2.668
1/16	64	1.15478	2.476
1/32	32	1.074607	2.3874
1/64	16	1.036633	2.3445
1/128	8	1.018152	2.3234 ²¹¹
1/256	4	1.0090350	2.3130 ¹⁰⁴
1/512	2	1.0045073	2.3077 ⁵³
1/1024	1	1.0022511	2.3051 ²⁶
			↓ 26
$\Delta/1024$ ($\Delta \rightarrow 0$)	Δ	$1 + .0022486\Delta$	2.3025

Now let us see how to calculate logarithms. We begin by computing successive square roots of 10, by cut and try. The results are shown in Table 22-1. The powers of 10 are given in the first column, and the result, 10^s , is given in the third column. Thus $10^1 = 10$. The one-half power of 10 we can easily work out, because that is the square root of 10, and there is a known, simple process for taking square roots of any number.* Using this process, we find the first square root to be 3.16228. What good is that? It already tells us something, it tells us how to take $10^{0.5}$, so we now know at least *one* logarithm, if we happen to need the logarithm of 3.16228, we know the answer is close to 0.50000. But we must do a little bit better than that; we clearly need more information. So we take the square root again, and find $10^{1/4}$, which is 1.77828. Now we have the logarithm of more numbers than we had before, 1.250 is the logarithm of 17.78 and, incidentally, if it happens that somebody asks for $10^{0.75}$, we can get it, because that is $10^{(0.5+0.25)}$; it is therefore the product of the second and third numbers. If we can get enough numbers in column s to be able to make up almost any number, then by multiplying the proper things in column 3, we can get 10 to any power; that is the plan. So we evaluate ten successive square roots of 10, and that is the main work which is involved in the calculations.

Why don't we keep on going for more and more accuracy? Because we begin to notice something. When we raise 10 to a very small power, we get 1 plus a small amount. The reason for this is clear, because we are going to have to take the

* There is a definite arithmetic procedure, but the easiest way to find the square root of any number N is to choose some a fairly close, find N/a , average $a' = \frac{1}{2}[a + (N/a)]$, and use this average a' for the next choice for a . The convergence is very rapid—the number of significant figures doubles each time.

1000th power of $10^{1/1000}$ to get back to 10, so we had better not start with too big a number; it has to be close to 1. What we notice is that the small numbers that are added to 1 begin to look as though we are merely dividing by 2 each time; we see 1815 becomes 903, then 450, 225; so it is clear that, to an excellent approximation, if we take another root, we shall get 1.00112 something, and rather than actually *take* all the square roots, we *guess* at the ultimate limit. When we take a small fraction Δ of 1024 as Δ approaches zero, what will the answer be? Of course it will be some number close to 0.0022511Δ . Not exactly 0.0022511Δ , however—we can get a better value by the following trick: we subtract the 1, and then divide by the power s . This ought to correct all the excesses to the same value. We see that they are very closely equal. At the top of the table they are not equal, but as they come down, they get closer and closer to a constant value. What is the value? Again we look to see how the series is going, how it has changed with s . It changed by 211, by 104, by 53, by 26. These changes are obviously half of each other, very closely, as we go down. Therefore, if we kept going, the changes would be 13, 7, 3, 2 and 1, more or less, or a total of 26. Thus we have only 26 more to go, and so we find that the true number is 2.3025. (Actually, we shall later see that the *exact* number should be 2.3026, but to keep it realistic, we shall not alter anything in the arithmetic.) From this table we can now calculate any power of 10, by compounding the power out of 1024ths.

Let us now actually calculate a logarithm, because the process we shall use is where logarithm tables actually come from. The procedure is shown in Table 22-2, and the numerical values are shown in Table 22-1 (columns 2 and 3).

Table 22-2

Calculation of a logarithm: $\log_{10} 2$

$$\begin{aligned}
 2 \div 1.77828 &= 1.124682 \\
 1.124682 \div 1.074607 &= 1.046598, \text{ etc.} \\
 \therefore 2 &= (1.77828)(1.074607)(1.036633)(1.090350)(1.000573) \\
 &= 10 \left[\frac{1}{1024} (256 + 32 + 16 + 4 + 0.254) \right] = 10 \left[\frac{308.254}{1024} \right] \\
 &= 10^{0.30103} \qquad \left(\frac{573}{2249} = 0.254 \right) \\
 \therefore \log_{10} 2 &= 0.30103
 \end{aligned}$$

Suppose we want the logarithm of 2. That is, we want to know to what power we must raise 10 to get 2. Can we raise 10 to the $1/2$ power? No; that is too big. In other words, we can see that the answer is going to be bigger than $1/4$, and less than $1/2$. Let us take the factor $10^{1/4}$ out; we divide 2 by 1.778 . . . , and get 1.124 . . . , and so on, and now we know that we have taken away 0.250000 from the logarithm. The number 1.124 . . . is now the number whose logarithm we need. When we are finished we shall add back the $1/4$, or 256/1024. Now we look in the table for the next number just below 1.124 . . . , and that is 1.074607. We therefore divide by 1.074607 and get 1.046598. From that we discover that 2 can be made up of a product of numbers that are in Table 22-1, as follows:

$$2 = (1.77828)(1.074607)(1.036633)(1.0090350)(1.000573).$$

There was one factor (1.000573) left over, naturally, which is beyond the range of our table. To get the logarithm of this factor, we use our result that $10^{\Delta/1024} \approx 1 + 2.3025 \Delta/1024$. We find $\Delta = 0.254$. Therefore our answer is 10 to the following power: $(256 + 32 + 16 + 4 + 0.254)/1024$. Adding those together, we get 308.254/1024. Dividing, we get 0.30103, so we know that the $\log_{10} 2 = 0.30103$, which happens to be right to 5 figures!

This is how logarithms were originally computed by Mr. Briggs of Halifax, in 1620. He said, "I computed successively 54 square roots of 10." We know he

really computed only the first 27, because the rest of them can be obtained by this trick with Δ . His work involved calculating the square root of 10 twenty-seven times, which is not much more than the ten times we did; however, it was more work because he calculated to sixteen decimal places, and then reduced his answer to fourteen when he published it, so that there were no rounding errors. He made tables of logarithms to fourteen decimal places by this method, which is quite tedious. But all logarithm tables for three hundred years were borrowed from Mr. Briggs' tables by reducing the number of decimal places. Only in modern times, with the WPA and computing machines, have new tables been independently computed. There are much more efficient methods of computing logarithms today, using certain series expansions.

In the above process, we discovered something rather interesting, and that is that for very small powers ϵ we can calculate 10^ϵ easily; we have discovered that $10^\epsilon = 1 + 2.3025\epsilon$, by sheer numerical analysis. Of course this also means that $10^{n/2.3025} = 1 + n$ if n is very small. Now logarithms to any other base are merely multiples of logarithms to the base 10. The base 10 was used only because we have 10 fingers, and the arithmetic of it is easy, but if we ask for a mathematically natural base, one that has nothing to do with the number of fingers on human beings, we might try to change our *scale* of logarithms in some convenient and natural manner, and the method which people have chosen is to redefine the logarithms by multiplying all the logarithms to the base 10 by 2.3025 . . . This then corresponds to using some other base, and this is called the *natural* base, or base e . Note that $\log_e(1 + n) \approx n$, or $e^n \approx 1 + n$ as $n \rightarrow 0$.

It is easy enough to find out what e is: $e = 10^{1/2.3025}$ or $10^{0.434294\dots}$, an irrational power. Our table of the successive square roots of 10 can be used to compute, not just logarithms, but also 10 to any power, so let us use it to calculate this natural base e . For convenience we transform $0.434294\dots$ into $444.73/1024$. Now, 444.73 is $256 + 128 + 32 + 16 + 2 + 0.73$. Therefore e , since it is an exponent of a sum, will be a product of the numbers

$$(1.77828)(1.33352)(1.074607)(1.036633)(1.018152)(1.009035)(1.001643) = 2.7184.$$

(The only problem is the last one, which is 0.73, and which is not in the table, but we know that if Δ is small enough, the answer is $1 + 2.3025 \Delta$.) When we multiply all these together, we get 2.7184 (it should be 2.7183, but it is good enough). The use of such tables, then, is the way in which irrational powers and the logarithms of irrational numbers are all calculated. That takes care of the irrationals.

22-5 Complex numbers

Now it turns out that after all that work we *still* cannot solve every equation! For instance, what is the square root of -1 ? Suppose we have to find $x^2 = -1$. The square of no rational, of no irrational, of *nothing* that we have discovered so far, is equal to -1 . So we again have to generalize our numbers to a still wider class. Let us suppose that a specific solution of $x^2 = -1$ is called something, we shall call it i ; i has the property, by definition, that its square is -1 . That is about all we are going to say about it; of course, there is more than one root of the equation $x^2 = -1$. Someone could write i , but another could say, "No, I prefer $-i$. My i is minus your i ." It is just as good a solution, and since the only definition that i has is that $i^2 = -1$, it must be true that any equation we can write is equally true if the sign of i is changed everywhere. This is called taking the *complex conjugate*. Now we are going to make up numbers by adding successive i 's, and multiplying i 's by numbers, and adding other numbers, and so on, according to all of our rules. In this way we find that numbers can all be written in the form $p + iq$, where p and q are what we call *real* numbers, i.e., the numbers we have been defining up until now. The number i is called the *unit imaginary* number. Any real multiple of i is called *pure imaginary*. The most general number, a , is of the form $p + iq$ and is called a *complex number*. Things do not get any worse if, for instance, we multiply two such numbers, let us say $(r + is)(p + iq)$. Then,

using the rules, we get

$$\begin{aligned}
 (r + is)(p + iq) &= rp + r(iq) + (is)p + (is)(iq) \\
 &= rp + i(rq) + i(sp) + (ii)(sq) \\
 &= (rp - sq) + i(rq + sp),
 \end{aligned} \tag{22.4}$$

since $ii = i^2 = -1$. Therefore all the numbers that now belong in the rules (22.1) have this mathematical form.

Now you say, “This can go on forever! We have defined powers of imaginaries and all the rest, and when we are all finished, somebody else will come along with another equation which cannot be solved, like $x^6 + 3x^2 = -2$. Then we have to generalize all over again!” But it turns out that *with this one more invention*, just the square root of -1 , *every algebraic equation can be solved!* This is a fantastic fact, which we must leave to the Mathematics Department to prove. The proofs are very beautiful and very interesting, but certainly not self-evident. In fact, the most obvious supposition is that we are going to have to invent again and again and again. But the greatest miracle of all is that we do not. This is the last invention. After this invention of complex numbers, we find that the rules still work with complex numbers, and we are finished inventing new things. We can find the complex power of any complex number, we can solve any equation that is written algebraically, in terms of a finite number of those symbols. We do not find any new numbers. The square root of i , for instance, has a definite result, it is not something new; and i^i is something. We will discuss that now.

We have already discussed multiplication, and addition is also easy; if we add two complex numbers, $(p + iq) + (r + is)$, the answer is $(p + r) + i(q + s)$. Now we can add and multiply complex numbers. But the real problem, of course, is to compute *complex powers of complex numbers*. It turns out that the problem is actually no more difficult than computing complex powers of real numbers. So let us concentrate now on the problem of calculating 10 to a complex power, not just an irrational power, but $10^{(r+is)}$. Of course, we must at all times use our rules (22.1) and (22.2). Thus

$$10^{(r+is)} = 10^r 10^{is}. \tag{22.5}$$

But 10^r we already know how to compute, and we can always multiply anything by anything else; therefore the problem is to compute only 10^{is} . Let us call it some complex number, $x + iy$. Problem: given s , find x , find y . Now if

$$10^{is} = x + iy,$$

then the complex conjugate of this equation must also be true, so that

$$10^{-is} = x - iy.$$

(Thus we see that we can deduce a number of things without actually computing anything, by using our rules.) We deduce another interesting thing by multiplying these together:

$$10^{is} 10^{-is} = 10^0 = 1 = (x + iy)(x - iy) = x^2 + y^2. \tag{22.6}$$

Thus if we find x , we have y also.

Now the problem is *how* to compute 10 to an imaginary power. What guide is there? We may work over our rules until we can go no further, but here is a reasonable guide: if we can compute it for any particular s , we can get it for all the rest. If we know 10^{is} for any one s and then we want it for twice that s , we can square the number, and so on. But how can we find 10^{is} for even one special value of s ? To do so we shall make one additional assumption, which is not quite in the category of all the other rules, but which leads to reasonable results and permits us to make progress: when the power is small, we shall suppose that the “law” $10^\epsilon = 1 + 2.3025\epsilon$ is right, as ϵ gets very small, not only for real ϵ , *but for complex ϵ as well*. Therefore, we begin with the supposition that this law is true in general, and that tells us that $10^{is} = 1 + 2.3025 \cdot is$, for $s \rightarrow 0$. So we assume that if s is very small, say one part in 1024, we have a rather good approximation to 10^{is} .

Now we make a table by which we can compute *all* the imaginary powers of 10, that is, compute x and y . It is done as follows. The first power we start with is the $1/1024$ power, which we presume is very nearly $1 + 2.3025i/1024$. Thus we start with

$$10^{i/1024} = 1.00000 + 0.0022486i, \quad (22.7)$$

and if we keep multiplying the number by itself, we can get to a higher imaginary power. In fact, we may just reverse the procedure we used in making our logarithm table, and calculate the square, 4th power, 8th power, etc., of (22.7), and thus build up the values shown in Table 22-3. We notice an interesting thing, that the x numbers are positive at first, but then swing negative. We shall look into that a little bit more in a moment. But first we may be curious to find for what number s the real part of 10^{is} is zero. The y -value would be i , and so we would have $10^{is} = i$, or $is = \log_{10} i$. As an example of how to use this table, just as we calculated $\log_{10} 2$ before, let us now use Table 22-3 to find $\log_{10} i$.

Which of the numbers in Table 22-3 do we have to multiply together to get a pure imaginary result? After a little trial and error, we discover that to reduce x the most, it is best to multiply "512" by "128." This gives $0.13056 + 0.99144i$. Then we discover that we should multiply this by a number whose imaginary part is about equal to the size of the real part we are trying to remove. Thus we choose "64" whose y -value is 0.14349, since that is closest to 0.13056. This then gives $-0.01350 + 0.99993i$. Now we have overshoot, and must *divide* by $0.99996 + 0.00900i$. How do we do that? By changing the sign of i and multiplying by $0.99996 - 0.00900i$ (which works if $x^2 + y^2 = 1$). Continuing in this way, we find that the entire power to which 10 must be raised to give i is $i(512 + 128 + 64 - 4 - 2 + 0.20)/1024$, or $698.20i/1024$. If we raise 10 to that power, we can get i . Therefore $\log_{10} i = 0.68226i$.

22-6 Imaginary exponents

To further investigate the subject of taking complex imaginary powers, let us look at the powers of 10 taking *successive powers*, not doubling the power each time, in order to follow Table 22-3 further and to see what happens to those minus signs. This is shown in Table 22-4, in which we take $10^{i/8}$, and just keep multiplying it. We see that x decreases, passes through zero, swings almost to -1 (if we could get in between $p = 10$ and $p = 11$ it would obviously swing to -1), and swings back. The y -value is going back and forth too.

In Fig. 22-1 the dots represent the numbers that appear in Table 22-4, and the lines are just drawn to help you visually. So we see that the numbers x and y oscillate; 10^{is} *repeats itself*, it is a *periodic* thing, and as such, it is easy enough to explain, because if a certain power is i , then the fourth power of that would be i^2 *squared*. It would be $+1$ again, and therefore, since $10^{0.68i}$ is equal to i , by taking the fourth power we discover that $10^{2.72i}$ is equal to $+1$. Therefore, if we wanted $10^{3.00i}$, for instance, we could write it as $10^{2.72i}$ times $10^{0.28i}$. In other words, it has a period, it repeats. Of course, we recognize what the curves look like! They look like the sine and cosine, and we shall call them, for a while, the algebraic sine and algebraic cosine. However, instead of using the base 10, we shall put them into our natural base, which only changes the horizontal scale; so we denote $2.3025s$ by t , and write $10^{is} = e^{it}$, where t is a real number. Now $e^{it} = x + iy$, and we shall write this as the algebraic cosine of t plus i times the algebraic sine of t . Thus

$$e^{it} = \cos t + i \sin t. \quad (22.8)$$

What are the properties of $\cos t$ and $\sin t$? First, we know, for instance, that $x^2 + y^2$ must be 1; we have proved that before, and it is just as true for base e as for base 10. Therefore $\cos^2 t + \sin^2 t = 1$. We also know that, for small t , $e^{it} = 1 + it$, and therefore $\cos t$ is nearly 1, and $\sin t$ is nearly t , and so it goes, that *all of the various properties of these remarkable functions*, which come from taking imaginary powers, *are the same as the sine and cosine of trigonometry*.

Table 22-3

Successive Squares of

$$10^{i/1024} = 1 + 0.0022486i$$

Power is	1024s	10^{is}
$i/1024$	1	$1.00000 + 0.00225i^*$
$i/512$	2	$1.00000 + 0.00450i$
$i/256$	4	$0.99996 + 0.00900i$
$i/128$	8	$0.99984 + 0.01800i$
$i/64$	16	$0.99936 + 0.03599i$
$i/32$	32	$0.99742 + 0.07193i$
$i/16$	64	$0.98967 + 0.14349i$
$i/8$	128	$0.95885 + 0.28402i$
$i/4$	256	$0.83872 + 0.54467i$
$i/2$	512	$0.40679 + 0.91365i$
$i/1$	1024	$-0.66928 + 0.74332i$

* Should be 0.0022486i

Table 22-4

Successive Powers of $10^{i/8}$

$p = \text{power} \cdot 8i$	$10^{ip/8}$
0	$1.00000 + 0.00000i$
1	$0.95882 + 0.28402i$
2	$0.83867 + 0.54465i$
3	$0.64944 + 0.76042i$
4	$0.40672 + 0.91356i$
5	$0.13050 + 0.99146i$
6	$-0.15647 + 0.98770i$
7	$-0.43055 + 0.90260i$
8	$-0.66917 + 0.74315i$
9	$-0.85268 + 0.52249i$
10	$-0.96596 + 0.25880i$
11	$-0.99969 - 0.02620i$
12	$-0.95104 - 0.30905i$
14	$-0.62928 - 0.77717i$
16	$-0.10447 - 0.99453i$
18	$+0.45454 - 0.89098i$
20	$+0.86648 - 0.49967i$
22	$+0.99884 + 0.05287i$
24	$+0.80890 + 0.58836i$

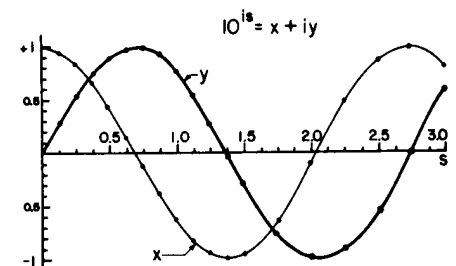


Figure 22-1

Is the period the same? Let us find out. e to what power is equal to i ? What is the logarithm of i to the base e ? We worked it out before, in the base 10 it was $0.68226i$, but when we change our logarithmic scale to e , we have to multiply by 2.3025, and if we do that it comes out 1.5709. So this will be called “algebraic $\pi/2$.” But, we see, it differs from the regular $\pi/2$ by only one place in the last point, and that, of course, is the result of errors in our arithmetic! So we have created two new functions in a purely algebraic manner, the cosine and the sine, which belong to algebra, and only to algebra. We wake up at the end to discover the very functions that are natural to geometry. So there is a connection, ultimately, between algebra and geometry.

We summarize with this, the most remarkable formula in mathematics:

$$e^{i\theta} = \cos \theta + i \sin \theta. \quad (22.9)$$

This is our jewel.

We may relate the geometry to the algebra by representing complex numbers in a plane; the horizontal position of a point is x , the vertical position of a point is y (Fig. 22-2). We represent every complex number, $x + iy$. Then if the radial distance to this point is called r and the angle is called θ , the algebraic law is that $x + iy$ is written in the form $re^{i\theta}$, where the geometrical relationships between x , y , r , and θ are as shown. This, then, is the unification of algebra and geometry.

When we began this chapter, armed only with the basic notions of integers and counting, we had little idea of the power of the processes of abstraction and generalization. Using the set of algebraic “laws,” or properties of numbers, Eq. (22.1), and the definitions of inverse operations (22.2), we have been able here, ourselves, to manufacture not only numbers but useful things like tables of logarithms, powers, and trigonometric functions (for these are what the imaginary powers of real numbers are), all merely by extracting ten successive square roots of ten!

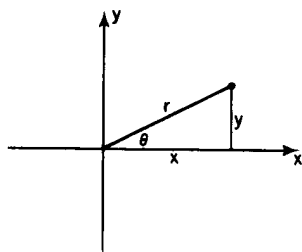


Fig. 22-2. $x + iy = re^{i\theta}$.

Resonance

23-1 Complex numbers and harmonic motion

In the present chapter we shall continue our discussion of the harmonic oscillator and, in particular, the forced harmonic oscillator, using a new technique in the analysis. In the preceding chapter we introduced the idea of complex numbers, which have real and imaginary parts and which can be represented on a diagram in which the ordinate represents the imaginary part and the abscissa represents the real part. If a is a complex number, we may write it as $a = a_r + ia_i$, where the subscript r means the real part of a , and the subscript i means the imaginary part of a . Referring to Fig. 23-1, we see that we may also write a complex number $a = x + iy$ in the form $x + iy = re^{i\theta}$, where $r^2 = x^2 + y^2 = (x + iy)(x - iy) = aa^*$. (The complex conjugate of a , written a^* , is obtained by reversing the sign of i in a .) So we shall represent a complex number in either of two forms, a real plus an imaginary part, or a magnitude r and a phase angle θ , so-called. Given r and θ , x and y are clearly $r \cos \theta$ and $r \sin \theta$ and, in reverse, given a complex number $x + iy$, $r = \sqrt{x^2 + y^2}$ and $\tan \theta = y/x$, the ratio of the imaginary to the real part.

We are going to apply complex numbers to our analysis of physical phenomena by the following trick. We have examples of things that oscillate; the oscillation may have a driving force which is a certain constant times $\cos \omega t$. Now such a force, $F = F_0 \cos \omega t$, can be written as the real part of a complex number $F = F_0 e^{i\omega t}$ because $e^{i\omega t} = \cos \omega t + i \sin \omega t$. The reason we do this is that it is easier to work with an exponential function than with a cosine. So the whole trick is to represent our oscillatory functions as the real parts of certain complex functions. The complex number F that we have so defined is not a real physical force, because no force in physics is really complex; actual forces have no imaginary part, only a real part. We shall, however, speak of the "force" $F_0 e^{i\omega t}$, but of course the actual force is the *real part* of that expression.

Let us take another example. Suppose we want to represent a force which is a cosine wave that is out of phase with a delayed phase Δ . This, of course, would be the real part of $F_0 e^{i(\omega t - \Delta)}$, but exponentials being what they are, we may write $e^{i(\omega t - \Delta)} = e^{i\omega t} e^{-i\Delta}$. Thus we see that the algebra of exponentials is much easier than that of sines and cosines; this is the reason we choose to use complex numbers. We shall often write

$$F = F_0 e^{-i\Delta} e^{i\omega t} = \hat{F} e^{i\omega t}. \quad (23.1)$$

We write a little caret (^) over the F to remind ourselves that this quantity is a complex number: here the number is

$$\hat{F} = F_0 e^{-i\Delta}.$$

Now let us solve an equation, using complex numbers, to see whether we can work out a problem for some real case. For example, let us try to solve

$$\frac{d^2 x}{dt^2} + \frac{kx}{m} = \frac{F}{m} = \frac{F_0}{m} \cos \omega t, \quad (23.2)$$

where F is the force which drives the oscillator and x is the displacement. Now, absurd though it may seem, let us suppose that x and F are actually complex numbers, for a mathematical purpose only. That is to say, x has a real part and an imaginary part times i , and F has a real part and an imaginary part times i .

23-1 Complex numbers and harmonic motion

23-2 The forced oscillator with damping

23-3 Electrical resonance

23-4 Resonance in nature

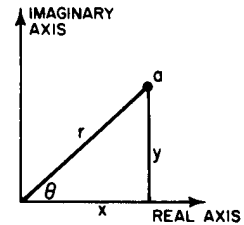


Fig. 23-1. A complex number may be represented by a point in the "complex plane."

Now if we had a solution of (23.2) with complex numbers, and substituted the complex numbers in the equation, we would get

$$\frac{d^2(x_r + ix_i)}{dt^2} + \frac{k(x_r + ix_i)}{m} = \frac{F_r + iF_i}{m}$$

or

$$\frac{d^2x_r}{dt^2} + \frac{kx_r}{m} + i\left(\frac{d^2x_i}{dt^2} + \frac{kx_i}{m}\right) = \frac{F_r}{m} + \frac{iF_i}{m}.$$

Now, since if two complex numbers are equal, their real parts must be equal *and* their complex parts must be equal, we deduce that *the real part of x satisfies the equation with the real part of the force*. We must emphasize, however, that this separation into a real part and an imaginary part is *not* valid in general, but is valid only for equations which are *linear*, that is, for equations in which x appears in every term only in the first power or the zeroth power. For instance, if there were in the equation a term λx^2 , then when we substitute $x_r + ix_i$, we would get $\lambda(x_r + ix_i)^2$, but when separated into real and imaginary parts this would yield $\lambda(x_r^2 - x_i^2)$ as the real part and $2i\lambda x_r x_i$ as the imaginary part. So we see that the real part of the equation would not involve just λx_r^2 , but also $-\lambda x_i^2$. In this case we get a different equation than the one we wanted to solve, with x_i , the completely artificial thing we introduced in our analysis, mixed in.

Let us now try our new method for the problem of the forced oscillator, that we already know how to solve. We want to solve Eq. (23.2) as before, but we say that we are going to solve

$$\frac{d^2x}{dt^2} + \frac{kx}{m} = \frac{\hat{F}e^{i\omega t}}{m}, \quad (23.3)$$

where $\hat{F}e^{i\omega t}$ is a complex number. Of course x will also be complex, but remember the rule: take the real part to find out what is really going on. So we try to solve (23.3) for the forced solution; we shall discuss other solutions later. The forced solution has the same frequency as the applied force, and has some amplitude of oscillation and some phase, and so it can be represented also by some complex number \hat{x} whose magnitude represents the swing of x and whose phase represents the time delay in the same way as for the force. Now a wonderful feature of an exponential function is that $d(\hat{x}e^{i\omega t})/dt = i\omega\hat{x}e^{i\omega t}$. When we differentiate an exponential function, we bring down the exponent as a simple multiplier. The second derivative does the same thing, it brings down another $i\omega$, and so it is very simple to write immediately, by inspection, what the equation is for \hat{x} : every time we see a differentiation, we simply multiply by $i\omega$. (Differentiation is now as easy as multiplication! This idea of using exponentials in linear differential equations is almost as great as the invention of logarithms, in which multiplication is replaced by addition. Here differentiation is replaced by multiplication.) Thus our equation becomes

$$(i\omega)^2\hat{x} + (k\hat{x}/m) = \hat{F}/m. \quad (23.4)$$

(We have cancelled the common factor $e^{i\omega t}$.) See how simple it is! Differential equations are immediately converted, by sight, into mere algebraic equations; we virtually have the solution by sight, that

$$\hat{x} = \frac{\hat{F}/m}{(k/m) - \omega^2},$$

since $(i\omega)^2 = -\omega^2$. This may be slightly simplified by substituting $k/m = \omega_0^2$, which gives

$$\hat{x} = \hat{F}/m(\omega_0^2 - \omega^2). \quad (23.5)$$

This, of course, is the solution we had before; for since $m(\omega_0^2 - \omega^2)$ is a real number, the phase angles of \hat{F} and of \hat{x} are the same (or perhaps 180° apart, if $\omega^2 > \omega_0^2$), as advertised previously. The magnitude of \hat{x} , which measures how far it oscillates, is related to the size of the \hat{F} by the factor $1/m(\omega_0^2 - \omega^2)$, and this factor becomes

enormous when ω is nearly equal to ω_0 . So we get a very strong response when we apply the right frequency ω (if we hold a pendulum on the end of a string and shake it at just the right frequency, we can make it swing very high).

23-2 The forced oscillator with damping

That, then, is how we analyze oscillatory motion with the more elegant mathematical technique. But the elegance of the technique is not at all exhibited in such a problem that can be solved easily by other methods. It is only exhibited when one applies it to more difficult problems. Let us therefore solve another, more difficult problem, which furthermore adds a relatively realistic feature to the previous one. Equation (23.5) tells us that if the frequency ω were exactly equal to ω_0 , we would have an infinite response. Actually, of course, no such infinite response occurs because some other things, like friction, which we have so far ignored, limits the response. Let us therefore add to Eq. (23.2) a friction term.

Ordinarily such a problem is very difficult because of the character and complexity of the frictional term. There are, however, many circumstances in which the frictional force is *proportional to the speed* with which the object moves. An example of such friction is the friction for slow motion of an object in oil or a thick liquid. There is no force when it is just standing still, but the faster it moves the faster the oil has to go past the object, and the greater is the resistance. So we shall assume that there is, in addition to the terms in (23.2), another term, a resistance force proportional to the velocity: $F_f = -c dx/dt$. It will be convenient, in our mathematical analysis, to write the constant c as m times γ to simplify the equation a little. This is just the same trick we use with k when we replace it by $m\omega_0^2$, just to simplify the algebra. Thus our equation will be

$$m(d^2x/dt^2) + c(dx/dt) + kx = F \quad (23.6)$$

or, writing $c = m\gamma$ and $k = m\omega_0^2$ and dividing out the mass m ,

$$(d^2x/dt^2) + \gamma(dx/dt) + \omega_0^2x = F/m. \quad (23.6a)$$

Now we have the equation in the most convenient form to solve. If γ is very small, that represents very little friction; if γ is very large, there is a tremendous amount of friction. How do we solve this new linear differential equation? Suppose that the driving force is equal to $F_0 \cos(\omega t + \Delta)$; we could put this into (23.6a) and try to solve it, but we shall instead solve it by our new method. Thus we write F as the real part of $\hat{F}e^{i\omega t}$ and x as the real part of $\hat{x}e^{i\omega t}$, and substitute these into Eq. (23.6a). It is not even necessary to do the actual substituting, for we can see by inspection that the equation would become

$$[(i\omega)^2\hat{x} + \gamma(i\omega)\hat{x} + \omega_0^2\hat{x}]e^{i\omega t} = (\hat{F}/m)e^{i\omega t}. \quad (23.7)$$

[As a matter of fact, if we tried to solve Eq. (23.6a) by our old straightforward way, we would really appreciate the magic of the “complex” method.] If we divide by $e^{i\omega t}$ on both sides, then we can obtain the response \hat{x} to the given force \hat{F} ; it is

$$\hat{x} = \hat{F}/m(\omega_0^2 - \omega^2 + i\gamma\omega). \quad (23.8)$$

Thus again \hat{x} is given by \hat{F} times a certain factor. There is no technical name for this factor, no particular letter for it, but we may call it R for discussion purposes:

$$R = \frac{1}{m(\omega_0^2 - \omega^2 + i\gamma\omega)}$$

and

$$\hat{x} = \hat{F}R. \quad (23.9)$$

(Although the letters γ and ω_0 are in very common use, this R has no particular name.) This factor R can either be written as $p + iq$, or as a certain magnitude ρ times $e^{i\theta}$. If it is written as a certain magnitude times $e^{i\theta}$, let us see what it means.

Now $\hat{F} = F_0 e^{i\Delta}$, and the actual force F is the real part of $F_0 e^{i\Delta} e^{i\omega t}$, that is, $F_0 \cos(\omega t + \Delta)$. Next, Eq. (23.9) tells us that \hat{x} is equal to $\hat{F}R$. So, writing $R = \rho e^{i\theta}$ as another name for R , we get

$$\hat{x} = R\hat{F} = \rho e^{i\theta} F_0 e^{i\Delta} = \rho F_0 e^{i(\theta + \Delta)}.$$

Finally, going even further back, we see that the physical x , which is the real part of the complex \hat{x} , is equal to the real part of $\rho F_0 e^{i(\theta + \Delta)} e^{i\omega t}$. But ρ and F_0 are real, and the real part of $e^{i(\theta + \Delta + \omega t)}$ is simply $\cos(\omega t + \Delta + \theta)$. Thus

$$x = \rho F_0 \cos(\omega t + \Delta + \theta). \quad (23.10)$$

This tells us that the amplitude of the response is the magnitude of the force F multiplied by a certain magnification factor, ρ ; this gives us the “amount” of oscillation. It also tells us, however, that x is not oscillating in phase with the force, which has the phase Δ , but is shifted by an extra amount θ . Therefore ρ and θ represent the size of the response and the phase shift of the response.

Now let us work out what ρ is. If we have a complex number, the square of the magnitude is equal to the number times its complex conjugate; thus

$$\begin{aligned} \rho^2 &= \frac{1}{m^2(\omega_0^2 - \omega^2 + i\gamma\omega)(\omega_0^2 - \omega^2 - i\gamma\omega)} \\ &= \frac{1}{m^2[(\omega^2 - \omega_0^2)^2 + \gamma^2\omega^2]}. \end{aligned} \quad (23.11)$$

In addition, the phase angle θ is easy to find, for if we write

$$1/R = 1/\rho e^{i\theta} = (1/\rho)e^{-i\theta} = m(\omega_0^2 - \omega^2 + i\gamma\omega),$$

we see that

$$\tan \theta = -\gamma\omega/(\omega_0^2 - \omega^2). \quad (23.12)$$

It is minus because $\tan(-\theta) = -\tan \theta$. A negative value for θ results for all ω , and this corresponds to the displacement x lagging the force F .

Figure 23-2 shows how ρ^2 varies as a function of frequency (ρ^2 is physically more interesting than ρ , because ρ^2 is proportional to the square of the amplitude, or more or less to the *energy* that is developed in the oscillator by the force). We see that if γ is very small, then $1/(\omega_0^2 - \omega^2)^2$ is the most important term, and the response tries to go up toward infinity when ω equals ω_0 . Now the “infinity” is not actually infinite because if $\omega = \omega_0$, then $1/\gamma^2\omega^2$ is still there. The phase shift varies as shown in Fig. 23-3.

In certain circumstances we get a slightly different formula than (23.8), also called a “resonance” formula, and one might think that it represents a different phenomenon, but it does not. The reason is that if γ is very small the most interesting part of the curve is near $\omega = \omega_0$, and we may replace (23.8) by an approximate formula which is very accurate if γ is small and ω is near ω_0 . Since $\omega_0^2 - \omega^2 = (\omega_0 - \omega)(\omega_0 + \omega)$, if ω is near ω_0 this is nearly the same as $2\omega_0(\omega_0 - \omega)$ and $\gamma\omega$ is nearly the same as $\gamma\omega_0$. Using these in (23.8), we see that $\omega_0^2 - \omega^2 + i\gamma\omega \approx 2\omega_0(\omega_0 - \omega + i\gamma/2)$, so that

$$\hat{x} \approx \hat{F}/2m\omega_0(\omega_0 - \omega + i\gamma/2) \quad \text{if } \gamma \ll \omega_0 \text{ and } \omega \approx \omega_0. \quad (23.13)$$

It is easy to find the corresponding formula for ρ^2 . It is

$$\rho^2 \approx 1/4m^2\omega_0^2[(\omega_0 - \omega)^2 + \gamma^2/4].$$

We shall leave it to the student to show the following: if we call the maximum height of the curve of ρ^2 vs. ω one unit, and we ask for the width $\Delta\omega$ of the curve, at one half the maximum height, the full width at half the maximum height of the curve is $\Delta\omega = \gamma$, supposing that γ is small. The resonance is sharper and sharper as the frictional effects are made smaller and smaller.

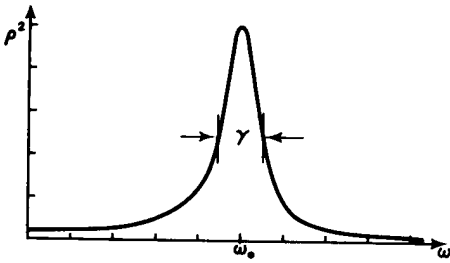


Fig. 23-2. Plot of ρ^2 versus ω .

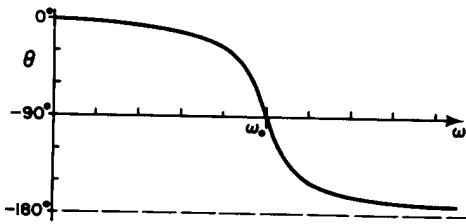


Fig. 23-3. Plot of θ versus ω .

As another measure of the width, some people use a quantity Q which is defined as $Q = \omega_0/\gamma$. The narrower the resonance, the higher the Q : $Q = 1000$ means a resonance whose width is only 1000th of the frequency scale. The Q of the resonance curve shown in Fig. 23-2 is 5.

The importance of the resonance phenomenon is that it occurs in many other circumstances, and so the rest of this chapter will describe some of these other circumstances.

23-3 Electrical resonance

The simplest and broadest technical applications of resonance are in electricity. In the electrical world there are a number of objects which can be connected to make electric circuits. These *passive circuit elements*, as they are often called, are of three main types, although each one has a little bit of the other two mixed in. Before describing them in greater detail, let us note that the whole idea of our mechanical oscillator being a mass on the end of a spring is only an approximation. All the mass is not actually at the “mass”; some of the mass is in the inertia of the spring. Similarly, all of the spring is not at the “spring”; the mass itself has a little elasticity, and although it may appear so, it is not *absolutely* rigid, and as it goes up and down, it flexes ever so slightly under the action of the spring pulling it. The same thing is true in electricity. There is an approximation in which we can lump things into “circuit elements” which are assumed to have pure, ideal characteristics. It is not the proper time to discuss that approximation here, we shall simply assume that it is true in the circumstances.

The three main kinds of circuit elements are the following. The first is called a *capacitor* (Fig. 23-4); an example is two plane metallic plates spaced a very small distance apart by an insulating material. When the plates are charged there is a certain voltage difference, that is, a certain difference in potential, between them. The same difference of potential appears between the terminals A and B , because if there were any difference along the connecting wire, electricity would flow right away. So there is a certain voltage difference V between the plates if there is a certain electric charge $+q$ and $-q$ on them, respectively. Between the plates there will be a certain electric field; we have even found a formula for it (Chapters 13 and 14):

$$V = \sigma d / \epsilon_0 = qd / \epsilon_0 A, \quad (23.14)$$

where d is the spacing and A is the area of the plates. Note that the potential difference is a linear function of the charge. If we do not have parallel plates, but insulated electrodes which are of any other shape, the difference in potential is still precisely proportional to the charge, but the constant of proportionality may not be so easy to compute. However, all we need to know is that the potential difference across a capacitor is *proportional to the charge*: $V = q/C$; the proportionality constant is $1/C$, where C is the *capacitance* of the object.

The second kind of circuit element is called a *resistor*; it offers resistance to the flow of electrical current. It turns out that metallic wires and many other substances resist the flow of electricity in this manner: if there is a voltage difference across a piece of some substance, there exists an electric current $I = dq/dt$ that is proportional to the electric voltage difference:

$$V = RI = R dq/dt. \quad (23.15)$$

The proportionality coefficient is called the *resistance* R . This relationship may already be familiar to you; it is Ohm's law.

If we think of the charge q on a capacitor as being analogous to the displacement x of a mechanical system, we see that the current, $I = dq/dt$, is analogous to velocity, $1/C$ is analogous to a spring constant k , and R is analogous to the resistive coefficient γ . Now it is very interesting that there exists another circuit element which is the analog of *mass*! This is a coil which builds up a magnetic field within itself when there is a current in it. A *changing* magnetic field develops in the coil a voltage that is proportional to dI/dt (this is how a transformer works, in

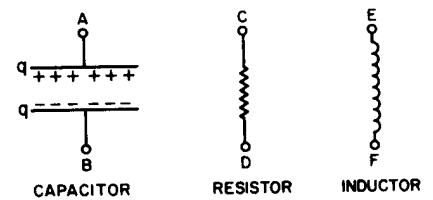


Fig. 23-4. The three passive circuit elements.

fact). The magnetic field is proportional to a current, and the induced voltage (so-called) in such a coil is proportional to the rate of change of the current:

$$V = L dI/dt = L d^2q/dt^2. \quad (23.16)$$

The coefficient L is the *self-inductance*, and is analogous to the mass in a mechanical oscillating circuit.

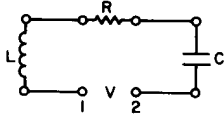


Fig. 23-5. An oscillatory electrical circuit with resistance, inductance, and capacitance.

Suppose we make a circuit in which we have connected the three circuit elements in series (Fig. 23-5); then the voltage across the whole thing from 1 to 2 is the work done in carrying a charge through, and it consists of the sum of several pieces: across the inductor, $V_L = L d^2q/dt^2$; across the resistance, $V_R = R dq/dt$; across the capacitor, $V_c = q/C$. The sum of these is equal to the applied voltage, V :

$$L d^2q/dt^2 + R dq/dt + q/C = V(t). \quad (23.17)$$

Now we see that this equation is exactly the same as the mechanical equation (23.6), and of course it can be solved in exactly the same manner. We suppose that $V(t)$ is oscillatory: we are driving the circuit with a generator with a pure sine wave oscillation. Then we can write our $V(t)$ as a complex \hat{V} with the understanding that it must be ultimately multiplied by $e^{i\omega t}$, and the real part taken in order to find the true V . Likewise, the charge q can thus be analyzed, and then in exactly the same manner as in Eq. (23.8) we write the corresponding equation: the second derivative of \hat{q} is $(i\omega)^2\hat{q}$; the first derivative is $(i\omega)\hat{q}$. Thus Eq. (23.17) translates to

$$\left[L(i\omega)^2 + R(i\omega) + \frac{1}{C} \right] \hat{q} = \hat{V}$$

or

$$\hat{q} = \frac{\hat{V}}{L(i\omega)^2 + R(i\omega) + \frac{1}{C}}$$

which we can write in the form

$$\hat{q} = \hat{V}/L(\omega_0^2 - \omega^2 + i\gamma\omega), \quad (23.18)$$

where $\omega_0^2 = 1/LC$ and $\gamma = R/L$. It is exactly the same denominator as we had in the mechanical case, with exactly the same resonance properties! The correspondence between the electrical and mechanical cases is outlined in Table 23-1.

Table 23-1

General characteristic	Mechanical property	Electrical property
indep. variable	time (t)	time (t)
dep. variable	position (x)	charge (q)
inertia	mass (m)	inductance (L)
resistance	drag coeff. ($c = \gamma m$)	resistance ($R = \gamma L$)
stiffness	stiffness (k)	(capacitance) $^{-1}$ ($1/C$)
resonant frequency	$\omega_0^2 = k/m$	$\omega_0^2 = 1/LC$
period	$t_0 = 2\pi\sqrt{m/k}$	$t_0 = 2\pi\sqrt{LC}$
figure of merit	$Q = \omega_0/\gamma$	$Q = \omega_0 L/R$

We must mention a small technical point. In the electrical literature, a different notation is used. (From one field to another, the subject is not really any different, but the way of writing the notations is often different.) First, j is commonly used instead of i in electrical engineering, to denote $\sqrt{-1}$. (After all, i must be the current!) Also, the engineers would rather have a relationship between \hat{V} and \hat{I} than between \hat{V} and \hat{q} , just because they are more used to it that way. Thus, since $\hat{I} = d\hat{q}/dt = i\omega\hat{q}$, we can just substitute $\hat{I}/i\omega$ for \hat{q} and get

$$\hat{V} = (i\omega L + R + 1/i\omega C)\hat{I} = \hat{Z}\hat{I}. \quad (23.19)$$

Another way is to rewrite Eq. (23.17), so that it looks more familiar; one often sees it written this way:

$$L \frac{dI}{dt} + RI + (1/C) \int^t I dt = V(t). \quad (23.20)$$

At any rate, we find the relation (23.19) between voltage \hat{V} and current \hat{I} which is just the same as (23.18) except divided by $i\omega$, and that produces Eq. (23.19). The quantity $R + i\omega L + 1/i\omega C$ is a complex number, and is used so much in electrical engineering that it has a name: it is called the *complex impedance*, \hat{Z} . Thus we can write $\hat{V} = \hat{Z}\hat{I}$. The reason that the engineers like to do this is that they learned something when they were young: $V = RI$ for resistances, when they only knew about resistances and DC. Now they have become more educated and have AC circuits, so they want the equation to look the same. Thus they write $\hat{V} = \hat{Z}\hat{I}$, the only difference being that the resistance is replaced by a more complicated thing, a complex quantity. So they insist that they cannot use what everyone else in the world uses for imaginary numbers, they have to use a j for that; it is a miracle that they did not insist also that the letter Z be an R ! (Then they get into trouble when they talk about current densities, for which they also use j . The difficulties of science are to a large extent the difficulties of notations, the units, and all the other artificialities which are invented by man, not by nature.)

23-4 Resonance in nature

Although we have discussed the electrical case in detail, we could also bring up case after case in many fields, and show exactly how the resonance equation is the same. There are many circumstances in nature in which something is "oscillating" and in which the resonance phenomenon occurs. We said that in an earlier chapter; let us now demonstrate it. If we walk around our study, pulling books off the shelves and simply looking through them to find an example of a curve that corresponds to Fig. 23-2 and comes from the same equation, what do we find? Just to demonstrate the wide range obtained by taking the smallest possible sample, it takes only five or six books to produce quite a series of phenomena which show resonances.

The first two are from mechanics, the first on a large scale: the atmosphere of the whole earth. If the atmosphere, which we suppose surrounds the earth evenly on all sides, is pulled to one side by the moon or, rather, squashed prolate into a double tide, and if we could then let it go, it would go sloshing up and down; it is an oscillator. This oscillator is *driven* by the moon, which is effectively revolving about the earth; any one component of the force, say in the x -direction, has a cosine component, and so the response of the earth's atmosphere to the tidal pull of the moon is that of an oscillator. The expected response of the atmosphere is shown in Fig. 23-6, curve b (curve a is another theoretical curve under discussion in the book from which this is taken out of context). Now one might think that we only have one point on this resonance curve, since we only have the one frequency, corresponding to the rotation of the earth under the moon, which occurs at a period of 12.42 hours—12 hours for the earth (the tide is a double bump), plus a little more because the moon is going around. But from the *size* of the atmospheric tides, and from the *phase*, the amount of delay, we can get both ρ and θ . From those we can get ω_0 and γ , and thus draw the entire curve! This is an example of very poor science. From two numbers we obtain two numbers, and from those two numbers we draw a beautiful curve, which of course goes through the very point that determined the curve! It is of no use *unless we can measure something else*, and in the case of geophysics that is often very difficult. But in this particular case there is another thing which we can show theoretically must have the same timing as the natural frequency ω_0 : that is, if someone disturbed the atmosphere, it would oscillate with the frequency ω_0 . Now there *was* such a sharp disturbance in 1883; the Krakatoa volcano exploded and half the island blew off, and it made such a terrific explosion in the atmosphere that the period of oscillation of the atmosphere could be measured. It came out to $10\frac{1}{2}$ hours. The ω_0 obtained from

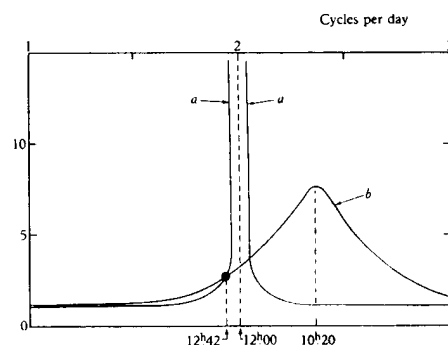


Fig. 23-6. Response of the atmosphere to external excitation. a is the required response if the atmospheric S_2 -tide is of gravitational origin; peak amplification is 100:1. b is derived from observed magnification and phase of M_2 -tide. [Munk and MacDonald, "Rotation of the Earth," Cambridge University Press (1960)]

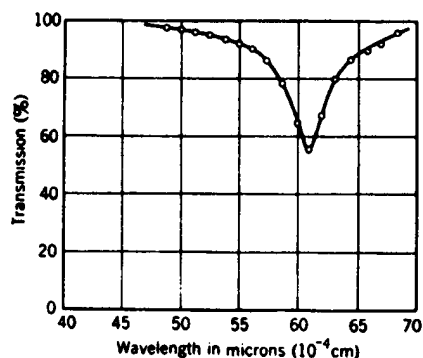


Fig. 23-7. Transmission of infrared radiation through a thin (0.17μ) sodium chloride film. [After R. B. Barnes, *Z. Physik* **75**, 723 (1932). Kittel, *Introduction to Solid State Physics*, Wiley, 1956.]

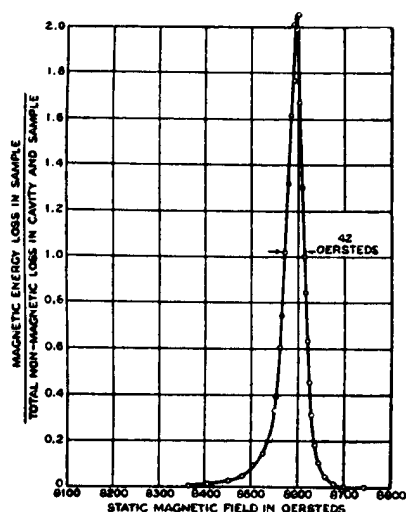


Fig. 23-8. Magnetic energy loss in paramagnetic organic compound as function of applied magnetic field intensity. [Holden et al., *Phys. Rev.* **75**, 1614 (1949)]

Fig. 23-6 comes out 10 hours and 20 minutes, so there we have at least one check on the reality of our understanding of the atmospheric tides.

Next we go to the small scale of mechanical oscillation. This time we take a sodium chloride crystal, which has sodium ions and chlorine ions next to each other, as we described in an early chapter. These ions are electrically charged, alternately plus and minus. Now there is an interesting oscillation possible. Suppose that we could drive all the plus charges to the right and all the negative charges to the left, and let go; they would then oscillate back and forth, the sodium lattice against the chlorine lattice. How can we ever drive such a thing? That is easy, for if we apply an electric field on the crystal, it will push the plus charge one way and the minus charge the other way! So, by having an external electric field we can perhaps get the crystal to oscillate. The frequency of the electric field needed is so high, however, that it corresponds to *infrared radiation*! So we try to find a resonance curve by measuring the absorption of infrared light by sodium chloride. Such a curve is shown in Fig. 23-7. The abscissa is not frequency, but is given in terms of wavelength, but that is just a technical matter, of course, since for a wave there is a definite relation between frequency and wavelength; so it is really a frequency scale, and a certain frequency corresponds to the resonant frequency.

But what about the width? What determines the width? There are many cases in which the width that is seen on the curve is not really the natural width γ that one would have theoretically. There are two reasons why there can be a wider curve than the theoretical curve. If the objects do not all have the same frequency, as might happen if the crystal were strained in certain regions, so that in those regions the oscillation frequency were slightly different than in other regions, then what we have is many resonance curves on top of each other; so we apparently get a wider curve. The other kind of width is simply this: perhaps we cannot measure the frequency precisely enough—if we open the slit of the spectrometer fairly wide, so although we thought we had only one frequency, we actually had a certain range $\Delta\omega$, then we may not have the resolving power needed to see a narrow curve. Offhand, we cannot say whether the width in Fig. 23-7 is natural, or whether it is due to inhomogeneities in the crystal or the finite width of the slit of the spectrometer.

Now we turn to a more esoteric example, and that is the swinging of a magnet. If we have a magnet, with north and south poles, in a constant magnetic field, the N end of the magnet will be pulled one way and the S end the other way, and there will in general be a torque on it, so it will vibrate about its equilibrium position, like a compass needle. However, the magnets we are talking about are *atoms*. These atoms have an angular momentum, the torque does not produce a simple motion in the direction of the field, but instead, of course, a *precession*. Now, looked at from the side, any one component is “swinging,” and we can disturb or drive that swinging and measure an absorption. The curve in Fig. 23-8 represents a typical such resonance curve. What has been done here is slightly different technically. The frequency of the lateral field that is used to drive this swinging is always kept the same, while we would have expected that the investigators would vary that and plot the curve. They could have done it that way, but technically it was easier for them to leave the frequency ω fixed, and change the strength of the constant magnetic field, which corresponds to changing ω_0 in our formula. They have plotted the resonance curve against ω_0 . Anyway, this is a typical resonance with a certain ω_0 and γ .

Now we go still further. Our next example has to do with atomic nuclei. The motions of protons and neutrons in nuclei are oscillatory in certain ways, and we can demonstrate this by the following experiment. We bombard a lithium atom with protons, and we discover that a certain reaction, producing γ -rays, actually has a very sharp maximum typical of resonance. We note in Fig. 23-9, however, one difference from other cases: the horizontal scale is not a frequency, it is an *energy*! The reason is that in quantum mechanics what we think of classically as the energy will turn out to be really related to a frequency of a wave amplitude. When we analyze something which in simple large-scale physics has to do with a frequency, we find that when we do quantum-mechanical experiments with atomic

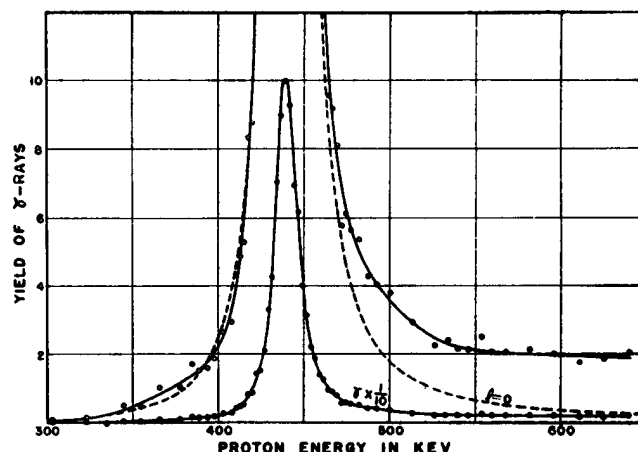


Fig. 23-9. The intensity of gamma-radiation from lithium as a function of the energy of the bombarding protons. The dashed curve is a theoretical one calculated for protons with an angular momentum $\ell = 0$. [Bonner and Evans, *Phys. Rev.* **73**, 666 (1948)]

matter, we get the corresponding curve as a function of energy. In fact, this curve is a demonstration of this relationship, in a sense. It shows that frequency and energy have some deep interrelationship, which of course they do.

Now we turn to another example which also involves a nuclear energy level, but now a much, much narrower one. The ω_0 in Fig. 23-10 corresponds to an energy of 100,000 electron volts, while the width γ is approximately 10^{-5} electron volt; in other words, this has a Q of 10^{10} ! When this curve was measured it was the largest Q of any oscillator that had ever been measured. It was measured by Dr. Mössbauer, and it was the basis of his Nobel prize. The horizontal scale here is velocity, because the technique for obtaining the slightly different frequencies was to use the Doppler effect, by moving the source relative to the absorber. One can see how delicate the experiment is when we realize that the speed involved is a few centimeters per second! On the actual scale of the figure, zero frequency would correspond to a point about 10^{10} cm to the left—slightly off the paper!

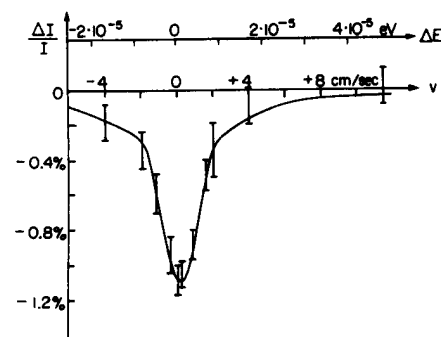


Fig. 23-10. [Courtesy of Dr. R. Mössbauer]

Finally, if we look in an issue of the *Physical Review*, say that of January 1, 1962, will we find a resonance curve? Every issue has a resonance curve, and Fig. 23-11 is the resonance curve for this one. This resonance curve turns out to be very interesting. It is the resonance found in a certain reaction among strange particles, a reaction in which a K^- and a proton interact. The resonance is detected by seeing how many of some kinds of particles come out, and depending on what and how many come out, one gets different curves, but of the same shape and with the peak at the same energy. We thus determine that there is a resonance at a certain energy for the K^- meson. That presumably means that there is some kind of a state, or condition, corresponding to this resonance, which can be attained by putting together a K^- and a proton. This is a new particle, or resonance. Today we do not know whether to call a bump like this a “particle” or simply a resonance. When there is a very *sharp* resonance, it corresponds to a very *definite* energy, just as though there were a particle of that energy present in nature. When the resonance gets wider, then we do not know whether to say there is a particle which does not last very long, or simply a resonance in the reaction probability. In the second chapter, this point is made about the particles, but when the second chapter was written this resonance was not known, so our chart should now have still another particle in it!

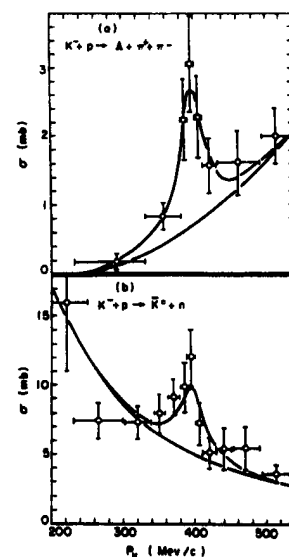


Fig. 23-11. Momentum dependence of the cross section for the reactions (a) $K^- + p \rightarrow \Lambda + \pi^+ + \pi^-$ and (b) $K^- + p \rightarrow K^0 + n$. The lower curves in (a) and (b) represent the presumed nonresonant backgrounds, while the upper curves contain in addition the superposed resonance. [Ferro-Luzzi *et al.*, *Phys. Rev., Lett.* **8**, 28 (1962)]

Transients

24-1 The energy of an oscillator

Although this chapter is entitled “transients,” certain parts of it are, in a way, part of the last chapter on forced oscillation. One of the features of a forced oscillation which we have not yet discussed is the *energy* in the oscillation. Let us now consider that energy.

In a mechanical oscillator, how much kinetic energy is there? It is proportional to the square of the velocity. Now we come to an important point. Consider an arbitrary quantity A , which may be the velocity or something else that we want to discuss. When we write $A = \hat{A}e^{i\omega t}$, a complex number, the true and honest A , in the physical world, is only the *real part*; therefore if, for some reason, we want to use the *square* of A , it is not right to square the complex number and then take the real part, because the real part of the square of a complex number is not just the square of the real part, but also involves the *imaginary* part. So when we wish to find the energy we have to get away from the complex notation for a while to see what the inner workings are.

Now the true physical A is the real part of $A_0 e^{i(\omega t + \Delta)}$, that is, $A = A_0 \cos(\omega t + \Delta)$, where \hat{A} , the complex number, is written as $A_0 e^{i\Delta}$. Now the square of this real physical quantity is $A^2 = A_0^2 \cos^2(\omega t + \Delta)$. The square of the quantity, then, goes up and down from a maximum to zero, like the square of the cosine. The square of the cosine has a maximum of 1 and a minimum of 0, and its average value is $1/2$.

In many circumstances we are not interested in the energy at any specific moment during the oscillation; for a large number of applications we merely want the average of A^2 , the *mean* of the square of A over a period of time large compared with the period of oscillation. In those circumstances, the average of the cosine squared may be used, so we have the following theorem: if A is represented by a complex number, then the mean of A^2 is equal to $\frac{1}{2}A_0^2$. Now A_0^2 is the square of the magnitude of the complex \hat{A} . (This can be written in many ways—some people like to write $|\hat{A}|^2$; others write, $\hat{A}\hat{A}^*$, \hat{A} times its complex conjugate.) We shall use this theorem several times.

Now let us consider the energy in a forced oscillator. The equation for the forced oscillator is

$$m d^2x/dt^2 + \gamma m dx/dt + m\omega_0^2 x = F(t). \quad (24.1)$$

In our problem, of course, $F(t)$ is a cosine function of t . Now let us analyze the situation: how much work is done by the outside force F ? The work done by the force per second, i.e., the power, is the force times the velocity. (We know that the differential work in a time dt is $F dx$, and the power is $F dx/dt$.) Thus

$$P = F \frac{dx}{dt} = m \left[\left(\frac{dx}{dt} \right) \left(\frac{d^2x}{dt^2} \right) + \omega_0^2 x \left(\frac{dx}{dt} \right) \right] + \gamma m \left(\frac{dx}{dt} \right)^2. \quad (24.2)$$

But the first two terms on the right can also be written as $d/dt[\frac{1}{2}m(dx/dt)^2 + \frac{1}{2}m\omega_0^2 x^2]$, as is immediately verified by differentiating. That is to say, the term in brackets is a pure derivative of two terms that are easy to understand—one is the kinetic energy of motion, and the other is the potential energy of the spring. Let us call this quantity the *stored energy*, that is, the energy stored in the oscillation. Suppose that we want the average power over many cycles when the oscillator is being forced and has been running for a long time. In the long run, the stored

24-1 The energy of an oscillator

24-2 Damped oscillations

24-3 Electrical transients

energy does not change—its derivative gives zero average effect. In other words, if we average the power in the long run, *all the energy ultimately ends up in the resistive term* $\gamma m(dx/dt)^2$. There is *some* energy stored in the oscillation, but that does not change with time, if we average over many cycles. Therefore the mean power $\langle P \rangle$ is

$$\langle P \rangle = \langle \gamma m(dx/dt)^2 \rangle. \quad (24.3)$$

Using our method of writing complex numbers, and our theorem that $\langle A^2 \rangle = \frac{1}{2}A_0^2$, we may find this mean power. Thus if $x = \hat{x}e^{i\omega t}$, then $dx/dt = i\omega\hat{x}e^{i\omega t}$. Therefore, in these circumstances, the average power could be written as

$$\langle P \rangle = \frac{1}{2}\gamma m\omega^2 x_0^2. \quad (24.4)$$

In the notation for electrical circuits, dx/dt is replaced by the current I (I is dq/dt , where q corresponds to x), and $m\gamma$ corresponds to the resistance R . Thus the rate of the energy loss—the power used up by the forcing function—is the resistance in the circuit times the average square of the current:

$$\langle P \rangle = R\langle I^2 \rangle = R \cdot \frac{1}{2}I_0^2. \quad (24.5)$$

This energy, of course, goes into heating the resistor; it is sometimes called the heating loss or the Joule heating.

Another interesting feature to discuss is how much energy is *stored*. That is not the same as the power, because although power was at first used to store up some energy, after that the system keeps on absorbing power, insofar as there are any heating (resistive) losses. At any moment there is a certain amount of stored energy, so we would like to calculate the mean stored energy $\langle E \rangle$ also. We have already calculated what the average of $(dx/dt)^2$ is, so we find

$$\begin{aligned} \langle E \rangle &= \frac{1}{2}m\langle (dx/dt)^2 \rangle + \frac{1}{2}m\omega_0^2\langle x^2 \rangle \\ &= \frac{1}{2}m(\omega^2 + \omega_0^2)\frac{1}{2}x_0^2. \end{aligned} \quad (24.6)$$

Now, when an oscillator is very efficient, and if ω is near ω_0 , so that $|\hat{x}|$ is large, the stored energy is very high—we can get a large stored energy from a relatively small force. The force does a great deal of work in getting the oscillation going, but then to keep it steady, all it has to do is to fight the friction. The oscillator can have a great deal of energy if the friction is very low, and even though it is oscillating strongly, not much energy is being lost. The efficiency of an oscillator can be measured by how much energy is stored, compared with how much work the force does per oscillation.

How does the stored energy compare with the amount of work that is done in one cycle? This is called the Q of the system, and Q is defined as 2π times the mean stored energy, divided by the work done per cycle. (If we were to say the work done per *radian* instead of per cycle, then the 2π disappears.)

$$Q = 2\pi \frac{\frac{1}{2}m(\omega^2 + \omega_0^2) \cdot \langle x^2 \rangle}{\gamma m\omega^2 \langle x^2 \rangle \cdot 2\pi/\omega} = \frac{\omega^2 + \omega_0^2}{2\gamma\omega}. \quad (24.7)$$

Q is not a very useful number unless it is very large. When it is relatively large, it gives a measure of how good the oscillator is. People have tried to define Q in the simplest and most useful way; various definitions differ a bit from one another, but if Q is very large, all definitions are in agreement. The most generally accepted definition is Eq. (24.7), which depends on ω . For a good oscillator, close to resonance, we can simplify (24.7) a little by setting $\omega = \omega_0$, and we then have $Q = \omega_0/\gamma$, which is the definition of Q that we used before.

What is Q for an electrical circuit? To find out, we merely have to translate L for m , R for $m\gamma$, and $1/C$ for $m\omega_0^2$ (see Table 23-1). The Q at resonance is $L\omega/R$, where ω is the resonance frequency. If we consider a circuit with a high Q , that means that the amount of energy stored in the oscillation is very large compared with the amount of work done per cycle by the machinery that drives the oscillations.

24-2 Damped oscillations

We now turn to our main topic of discussion: transients. By a *transient* is meant a solution of the differential equation when there is no force present, but when the system is not simply at rest. (Of course, if it is standing still at the origin with no force acting, that is a nice problem—it stays there!) Suppose the oscillation starts another way: say it was driven by a force for a while, and then we turn off the force. What happens then? Let us first get a rough idea of what will happen for a very high Q system. So long as a force is acting, the stored energy stays the same, and there is a certain amount of work done to maintain it. Now suppose we turn off the force, and no more work is being done; then the losses which are eating up the energy of the supply are no longer eating up its energy—there *is* no more driver. The losses will have to consume, so to speak, the energy that is stored. Let us suppose that $Q/2\pi = 1000$. Then the work done per cycle is $1/1000$ of the stored energy. Is it not reasonable, since it is oscillating with no driving force, that in one cycle the system will still lose a thousandth of its energy E , which ordinarily would have been supplied from the outside, and that it will continue oscillating, always losing $1/1000$ of its energy per cycle? So, as a guess, for a relatively high Q system, we would suppose that the following equation might be roughly right (we will later do it exactly, and it will turn out that it *was* right!):

$$dE/dt = -\omega E/Q. \quad (24.8)$$

This is rough because it is true only for large Q . In each radian the system loses a fraction $1/Q$ of the stored energy E . Thus in a given amount of time dt the energy will change by an amount $\omega dt/Q$, since the number of radians associated with the time dt is ωdt . What is the frequency? Let us suppose that the system moves so nicely, with hardly any force, that if we let go it will oscillate at essentially the same frequency all by itself. So we will guess that ω is the resonant frequency ω_0 . Then we deduce from Eq. (24.8) that the stored energy will vary as

$$E = E_0 e^{-\omega_0 t/Q} = E_0 e^{-\gamma t}. \quad (24.9)$$

This would be the measure of the *energy* at any moment. What would the formula be, roughly, for the amplitude of the oscillation as a function of the time? The same? No! The amount of energy in a spring, say, goes as the *square* of the displacement; the kinetic energy goes as the *square* of the velocity; so the total energy goes as the *square* of the displacement. Thus the displacement, the amplitude of oscillation, will decrease half as fast because of the square. In other words, we guess that the solution for the damped transient motion will be an oscillation of frequency close to the resonance frequency ω_0 , in which the amplitude of the sine-wave motion will diminish as $e^{-\gamma t/2}$:

$$x = A_0 e^{-\gamma t/2} \cos \omega_0 t. \quad (24.10)$$

This equation and Fig. 24-1 give us an idea of what we should expect; now let us try to analyze the motion *precisely* by solving the differential equation of the motion itself.

So, starting with Eq. (24.1), with no outside force, how do we solve it? Being physicists, we do not have to worry about the *method* as much as we do about what the solution *is*. Armed with our previous experience, let us try as a solution an exponential curve, $x = Ae^{i\alpha t}$. (Why do we try this? It is the easiest thing to differentiate!) We put this into (24.1) (with $F(t) = 0$), using the rule that each time we differentiate x with respect to time, we multiply by $i\alpha$. So it is really quite simple to substitute. Thus our equation looks like this:

$$(-\alpha^2 + i\gamma\alpha + \omega_0^2)Ae^{i\alpha t} = 0. \quad (24.11)$$

The net result must be zero for *all times*, which is impossible unless (a) $A = 0$, which is no solution at all—it stands still, or (b)

$$-\alpha^2 + i\alpha\gamma + \omega_0^2 = 0. \quad (24.12)$$

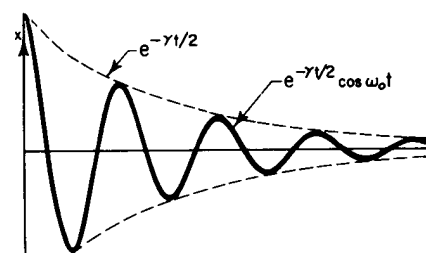


Fig. 24-1. A damped cosine oscillation.

If we can solve this and find an α , then we will have a solution in which A need not be zero!

$$\alpha = i\gamma/2 \pm \sqrt{\omega_0^2 - \gamma^2/4}. \quad (24.13)$$

For a while we shall assume that γ is fairly small compared with ω_0 , so that $\omega_0^2 - \gamma^2/4$ is definitely positive, and there is nothing the matter with taking the square root. The only bothersome thing is that we get *two* solutions! Thus

$$\alpha_1 = i\gamma/2 + \sqrt{\omega_0^2 - \gamma^2/4} = i\gamma/2 + \omega_\gamma \quad (24.14)$$

and

$$\alpha_2 = i\gamma/2 - \sqrt{\omega_0^2 - \gamma^2/4} = i\gamma/2 - \omega_\gamma. \quad (24.15)$$

Let us consider the first one, supposing that we had not noticed that the square root has two possible values. Then we know that a solution for x is $x_1 = Ae^{i\alpha_1 t}$, where A is any constant whatever. Now, in substituting α_1 , because it is going to come so many times and it takes so long to write, we shall call $\sqrt{\omega_0^2 - \gamma^2/4} = \omega_\gamma$. Thus $i\alpha_1 = -\gamma/2 + i\omega_\gamma$, and we get $x = Ae^{(-\gamma/2 + i\omega_\gamma)t}$, or what is the same, because of the wonderful properties of an exponential,

$$x_1 = Ae^{-\gamma t/2} e^{i\omega_\gamma t}. \quad (24.16)$$

First, we recognize this as an oscillation, an oscillation at a frequency ω_γ which is not *exactly* the frequency ω_0 , but is rather close to ω_0 if it is a good system. Second, the amplitude of the oscillation is decreasing exponentially! If we take, for instance, the real part of (24.16), we get

$$x_1 = Ae^{-\gamma t/2} \cos \omega_\gamma t. \quad (24.17)$$

This is very much like our guessed-at solution (24.10), except that the frequency really is ω_γ . This is the only error, so it is the same thing—we have the right idea. But everything is *not* all right! What is not all right is that *there is another solution*.

The other solution is α_2 , and we see that the difference is only that the sign of ω_γ is reversed:

$$x_2 = Be^{-\gamma t/2} e^{-i\omega_\gamma t}. \quad (24.18)$$

What does this mean? We shall soon prove that if x_1 and x_2 are each a possible solution of Eq. (24.1) with $F = 0$, then $x_1 + x_2$ is also a solution of the same equation! So the general solution x is of the mathematical form

$$x = e^{-\gamma t/2} (Ae^{i\omega_\gamma t} + Be^{-i\omega_\gamma t}). \quad (24.19)$$

Now we may wonder why we bother to give this other solution, since we were so happy with the first one all by itself. What is the extra one for, because of course we know we should only take the real part? *We* know that we must take the real part, but how did the *mathematics* know that we only wanted the real part? When we had a nonzero driving force $F(t)$, we put in an *artificial* force to go with it, and the *imaginary* part of the equation, so to speak, was driven in a definite way. But when we put $F(t) \equiv 0$, our convention that x should be only the real part of whatever we write down is purely our own, and the mathematical equations do not know it yet. The physical world *has* a real solution, but the answer that we were so happy with before is not real, it is *complex*. The *equation* does not know that we are arbitrarily going to take the real part, so it has to present us, so to speak, with a complex conjugate type of solution, so that by putting them together we can *make a truly real* solution; that is what α_2 is doing for us. In order for x to be real, $Be^{-i\omega_\gamma t}$ will have to be the complex conjugate of $Ae^{i\omega_\gamma t}$, so that the imaginary parts disappear. So it turns out that B is the complex conjugate of A , and our real solution is

$$x = e^{-\gamma t/2} (Ae^{i\omega_\gamma t} + A^* e^{-i\omega_\gamma t}). \quad (24.20)$$

So our real solution is an oscillation with a *phase shift* and a damping—just as advertised.

24-3 Electrical transients

Now let us see if the above really works. We construct the electrical circuit shown in Fig. 24-2, in which we apply to an oscilloscope the voltage across the inductance L after we suddenly turn on a voltage by closing the switch S . It is an oscillatory circuit, and it generates a transient of some kind. It corresponds to a circumstance in which we suddenly apply a force and the system starts to oscillate. It is the electrical analog of a damped mechanical oscillator, and we watch the oscillation on an oscilloscope, where we should see the curves that we were trying to analyze. (The horizontal motion of the oscilloscope is driven at a uniform speed, while the vertical motion is the voltage across the inductor. The rest of the circuit is only a technical detail. We would like to repeat the experiment many, many times, since the persistence of vision is not good enough to see only one trace on the screen. So we do the experiment again and again by closing the switch 60 times a second; each time we close the switch, we also start the oscilloscope horizontal sweep, and it draws the curve over and over.) In Figs. 24-3 to 24-6 we see examples of damped oscillations, actually photographed on an oscilloscope screen. Figure 24-3 shows a damped oscillation in a circuit which has a high Q , a small γ . It does not die out very fast; it oscillates many times on the way down.

But let us see what happens as we decrease Q , so that the oscillation dies out more rapidly. We can decrease Q by increasing the resistance R in the circuit. When we increase the resistance in the circuit, it dies out faster (Fig. 24-4). Then if we increase the resistance in the circuit still more, it dies out faster still (Fig. 24-5). But when we put in more than a certain amount, we cannot see any oscillation at all! The question is, is this because our eyes are not good enough? If we increase the resistance still more, we get a curve like that of Fig. 24-6, which does not appear to have any oscillations, except perhaps one. Now, how can we explain that by mathematics?

The resistance is, of course, proportional to the γ term in the mechanical device. Specifically, γ is R/L . Now if we increase the γ in the solutions (24.14) and (24.15) that we were so happy with before, chaos sets in when $\gamma/2$ exceeds ω_0 ; we must write it a different way, as

$$i\gamma/2 + i\sqrt{\gamma^2/4 - \omega_0^2} \quad \text{and} \quad i\gamma/2 - i\sqrt{\gamma^2/4 - \omega_0^2}.$$

Those are now the two solutions and, following the same line of mathematical reasoning as previously, we again find two solutions: $e^{i\alpha_1 t}$ and $e^{i\alpha_2 t}$. If we now substitute for α_1 , we get

$$x = Ae^{-(\gamma/2 + \sqrt{\gamma^2/4 - \omega_0^2})t},$$

a nice exponential decay with no oscillations. Likewise, the other solution is

$$x = Be^{-(\gamma/2 - \sqrt{\gamma^2/4 - \omega_0^2})t}.$$

Note that the square root cannot exceed $\gamma/2$, because even if $\omega_0 = 0$, one term just equals the other. But ω_0^2 is taken away from $\gamma^2/4$, so the square root is less than $\gamma/2$, and the term in parentheses is, therefore, always a positive number. Thank goodness! Why? Because if it were negative, we would find e raised to a *positive* factor times t , which would mean it was exploding! In putting more and more resistance into the circuit, we know it is not going to explode—quite the contrary. So now we have two solutions, each one by itself a dying exponential, but one having a much faster “dying rate” than the other. The general solution is of course a combination of the two; the coefficients in the combination depending upon how the motion starts—what the initial conditions of the problem are. In the particular way this circuit happens to be starting, the A is negative and the B is positive, so we get the difference of two exponential curves.

Now let us discuss how we can find the two coefficients A and B (or A and A^*), if we know how the motion was started.

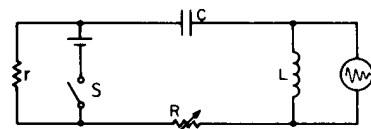


Fig. 24-2. An electrical circuit for demonstrating transients.

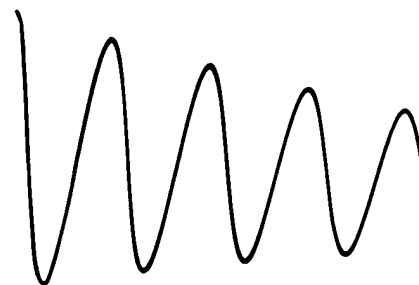


Figure 24-3

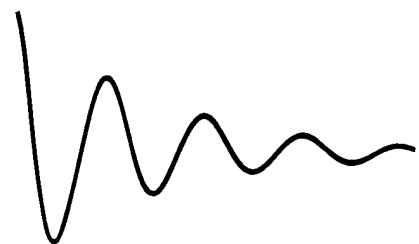


Figure 24-4

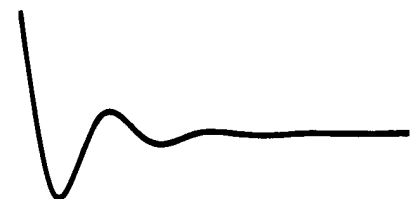


Figure 24-5



Figure 24-6

Suppose that at $t = 0$ we know that $x = x_0$, and $dx/dt = v_0$. If we put $t = 0$, $x = x_0$, and $dx/dt = v_0$ into the expressions

$$x = e^{-\gamma t/2}(Ae^{i\omega_\gamma t} + A^*e^{-i\omega_\gamma t}),$$

$$dx/dt = e^{-\gamma t/2}[(-\gamma/2 + i\omega_\gamma)Ae^{i\omega_\gamma t} + (-\gamma/2 - i\omega_\gamma)A^*e^{-i\omega_\gamma t}],$$

we find, since $e^0 = e^{i0} = 1$,

$$x_0 = A + A^* = 2A_R,$$

$$v_0 = (-\gamma/2)(A + A^*) + i\omega_\gamma(A - A^*)$$

$$= -\gamma x_0/2 + i\omega_\gamma(2iA_I),$$

where $A = A_R + iA_I$, and $A^* = A_R - iA_I$. Thus we find

$$A_R = x_0/2$$

and

$$A_I = (v_0 + \gamma x_0/2)/2\omega_\gamma. \quad (24.21)$$

This completely determines A and A^* , and therefore the complete curve of the transient solution, in terms of how it begins. Incidentally, we can write the solution another way if we note that

$$e^{i\theta} + e^{-i\theta} = 2 \cos \theta \quad \text{and} \quad e^{i\theta} - e^{-i\theta} = 2i \sin \theta.$$

We may then write the complete solution as

$$x = e^{-\gamma t/2} \left[x_0 \cos \omega_\gamma t + \frac{v_0 + \gamma x_0/2}{\omega_\gamma} \sin \omega_\gamma t \right], \quad (24.22)$$

where $\omega_\gamma = +\sqrt{\omega_0^2 - \gamma^2/4}$. This is the mathematical expression for the way an oscillation dies out. We shall not make direct use of it, but there are a number of points we should like to emphasize that are true in more general cases.

First of all the behavior of such a system with no external force is expressed by a sum, or superposition, of pure exponentials in time (which we wrote as $e^{i\alpha t}$). This is a good solution to try in such circumstances. The values of α may be complex in general, the imaginary parts representing damping. Finally the intimate mathematical relation of the sinusoidal and exponential function discussed in Chapter 22 often appears physically as a change from oscillatory to exponential behavior when some physical parameter (in this case resistance, γ) exceeds some critical value.

Linear Systems and Review

25-1 Linear differential equations

In this chapter we shall discuss certain aspects of oscillating systems that are found somewhat more generally than just in the particular systems we have been discussing. For our particular system, the differential equation that we have been solving is

$$m \frac{d^2x}{dt^2} + \gamma m \frac{dx}{dt} + m\omega_0^2 x = F(t). \quad (25.1)$$

Now this particular combination of “operations” on the variable x has the interesting property that if we substitute $(x + y)$ for x , then we get the sum of the same operations on x and y ; or, if we multiply x by a , then we get just a times the same combination. This is easy to prove. Just as a “shorthand” notation, because we get tired of writing down all those letters in (25.1), we shall use the symbol $\underline{L}(x)$ instead. When we see this, it means the left-hand side of (25.1), with x substituted in. With this system of writing, $\underline{L}(x + y)$ would mean the following:

$$\underline{L}(x + y) = m \frac{d^2(x + y)}{dt^2} + \gamma m \frac{d(x + y)}{dt} + m\omega_0^2(x + y). \quad (25.2)$$

(We underline the \underline{L} so as to remind ourselves that it is not an ordinary function.) We sometimes call this an *operator notation*, but it makes no difference what we call it, it is just “shorthand.”

Our first statement was that

$$\underline{L}(x + y) = \underline{L}(x) + \underline{L}(y), \quad (25.3)$$

which of course follows from the fact that $a(x + y) = ax + ay$, $d(x + y)/dt = dx/dt + dy/dt$, etc.

Our second statement was, for constant a ,

$$\underline{L}(ax) = a\underline{L}(x). \quad (25.4)$$

[Actually, (25.3) and (25.4) are very closely related, because if we put $x + x$ into (25.3), this is the same as setting $a = 2$ in (25.4), and so on.]

In more complicated problems, there may be more derivatives, and more terms in \underline{L} ; the question of interest is whether the two equations (25.3) and (25.4) are maintained or not. If they are, we call such a problem a *linear* problem. In this chapter we shall discuss some of the properties that exist because the system is linear, to appreciate the generality of some of the results that we have obtained in our special analysis of a special equation.

Now let us study some of the properties of linear differential equations, having illustrated them already with the specific equation (25.1) that we have studied so closely. The first property of interest is this: suppose that we have to solve the differential equation for a transient, the free oscillation with no driving force. That is, we want to solve

$$\underline{L}(x) = 0. \quad (25.5)$$

Suppose that, by some hook or crook, we have found a particular solution, which we shall call x_1 . That is, we have an x_1 for which $\underline{L}(x_1) = 0$. Now we notice that ax_1 is also a solution to the same equation; we can multiply this special solution by any constant whatever, and get a new solution. In other words, if we had a

25-1 Linear differential equations

25-2 Superposition of solutions

25-3 Oscillations in linear systems

25-4 Analogs in physics

25-5 Series and parallel impedances

motion of a certain “size,” then a motion twice as “big” is again a solution. *Proof:* $\underline{L}(ax_1) = a\underline{L}(x_1) = a \cdot 0 = 0$.

Next, suppose that, by hook or by crook, we have not only found *one* solution x_1 , but also another solution, x_2 . (Remember that when we substituted $x = e^{iat}$ for finding the transients, we found *two* values for α , that is, two solutions, x_1 and x_2 .) Now let us show that the combination $(x_1 + x_2)$ is also a solution. In other words, if we put $x = x_1 + x_2$, x is again a solution of the equation. Why? Because, if $\underline{L}(x_1) = 0$ and $\underline{L}(x_2) = 0$, then $\underline{L}(x_1 + x_2) = \underline{L}(x_1) + \underline{L}(x_2) = 0 + 0 = 0$. So if we have found a number of solutions for the motion of a linear system we can add them together.

Combining these two ideas, we see, of course, that we can also add six of one and two of the other: if x_1 is a solution, so is αx_1 . Therefore any sum of these two solutions, such as $(\alpha x_1 + \beta x_2)$, is also a solution. If we happen to be able to find three solutions, then we find that any combination of the three solutions is again a solution, and so on. It turns out that the number of what we call *independent solutions** that we have obtained for our oscillator problem is only *two*. The number of independent solutions that one finds in the general case depends upon what is called the number of *degrees of freedom*. We shall not discuss this in detail now, but if we have a second-order differential equation, there are only two independent solutions, and we have found both of them; so we have the most general solution.

Now let us go on to another proposition, which applies to the situation in which the system is subjected to an outside force. Suppose the equation is

$$\underline{L}(x) = F(t), \quad (25.6)$$

and suppose that we have found a special solution of it. Let us say that Joe’s solution is x_J , and that $\underline{L}(x_J) = F(t)$. Suppose we want to find yet another solution; suppose we add to Joe’s solution one of those that was a solution of the free equation (25.5), say x_1 . Then we see by (25.3) that

$$\underline{L}(x_J + x_1) = \underline{L}(x_J) + \underline{L}(x_1) = F(t) + 0 = F(t). \quad (25.7)$$

Therefore, to the “forced” solution we can add any “free” solution, and we still have a solution. The free solution is called a *transient* solution.

When we have no force acting, and suddenly turn one on, we do not immediately get the steady solution that we solved for with the sine wave solution, but for a while there is a transient which sooner or later dies out, if we wait long enough. The “forced” solution does not die out, since it keeps on being driven by the force. Ultimately, for long periods of time, the solution is unique, but initially the motions are different for different circumstances, depending on how the system was started.

25-2 Superposition of solutions

Now we come to another interesting proposition. Suppose that we have a certain particular driving force F_a (let us say an oscillatory one with a certain $\omega = \omega_a$, but our conclusions will be true for any functional form of F_a) and we have solved for the forced motion (with or without the transients; it makes no difference). Now suppose some other force is acting, let us say F_b , and we solve the same problem, but for this different force. Then suppose someone comes along and says, “I have a new problem for you to solve; I have the force $F_a + F_b$.” Can we do it? Of course we can do it, because the solution is the sum of the two solutions x_a and x_b for the forces taken separately—a most remarkable circumstance indeed. If we use (25.3), we see that

$$\underline{L}(x_a + x_b) = \underline{L}(x_a) + \underline{L}(x_b) = F_a(t) + F_b(t). \quad (25.8)$$

* Solutions which cannot be expressed as linear combinations of each other are called independent.

This is an example of what is called the *principle of superposition* for linear systems, and it is very important. It means the following: if we have a complicated force which can be broken up in any convenient manner into a sum of separate pieces, each of which is in some way simple, in the sense that for each special piece into which we have divided the force we can solve the equation, then the answer is available for the *whole* force, because we may simply add the pieces of the *solution* back together, in the same manner as the total *force* is compounded out of pieces (Fig. 25-1).

Let us give another example of the principle of superposition. In Chapter 12 we said that it was one of the great facts of the laws of electricity that if we have a certain distribution of charges q_a and calculate the electric field E_a arising from these charges at a certain place P , and if, on the other hand, we have another set of charges q_b and we calculate the field E_b due to these at the corresponding place, then if both charge distributions are present at the same time, the field E at P is the *sum* of E_a due to one set plus E_b due to the other. In other words, if we know the field due to a certain charge, then the field due to many charges is merely the vector sum of the fields of these charges taken individually. This is exactly analogous to the above proposition that if we know the result of two given forces taken at one time, then if the force is considered as a sum of them, the response is a sum of the corresponding individual responses.

The reason why this is true in electricity is that the great laws of electricity, Maxwell's equations, which determine the electric field, turn out to be differential equations which are *linear*, i.e., which have the property (25.3). What corresponds to the force is the *charge* generating the electric field, and the equation which determines the electric field in terms of the charge is linear.

As another interesting example of this proposition, let us ask how it is possible to "tune in" to a particular radio station at the same time as all the radio stations are broadcasting. The radio station transmits, fundamentally, an oscillating electric field of very high frequency which acts on our radio antenna. It is true that the amplitude of the oscillation of the field is changed, modulated, to carry the signal of the voice, but that is very slow, and we are not going to worry about it. When one hears "This station is broadcasting at a frequency of 780 kilocycles," this indicates that 780,000 oscillations per second is the frequency of the electric field of the station antenna, and this drives the electrons up and down at that frequency in our antenna. Now at the same time we may have another radio station in the same town radiating at a different frequency, say 550 kilocycles per second; then the electrons in our antenna are also being driven by that frequency. Now the question is, how is it that we can separate the signals coming into the one radio at 780 kilocycles from those coming in at 550 kilocycles? We certainly do not hear both stations at the same time.

By the principle of superposition, the response of the electric circuit in the radio, the first part of which is a linear circuit, to the forces that are acting due to the electric field $F_a + F_b$, is $x_a + x_b$. It therefore looks as though we will never disentangle them. In fact, the very proposition of superposition seems to insist that we cannot *avoid* having both of them in our system. But remember, for a *resonant* circuit, the response curve, the amount of x per unit F , as a function of the frequency, looks like Fig. 25-3. If it were a very high Q circuit, the response would show a very sharp maximum. Now suppose that the two stations are comparable in strength, that is, the two *forces* are of the same magnitude. The *response* that we get is the sum of x_a and x_b . But, in Fig. 25-3, x_a is tremendous, while x_b is small. So, in spite of the fact that the two signals are equal in strength, when they go through the sharp resonant circuit of the radio tuned for ω_a , the frequency of the transmission of one station, then the response to this station is much greater than to the other. Therefore the complete response, with both signals acting, is almost all made up of ω_a , and we have selected the station we want.

Now what about the tuning? How do we tune it? We change ω_0 by changing the L or the C of the circuit, because the frequency of the circuit has to do with the combination of L and C . In particular, most radios are built so that one can change the capacitance. When we retune the radio, we can make a new setting of

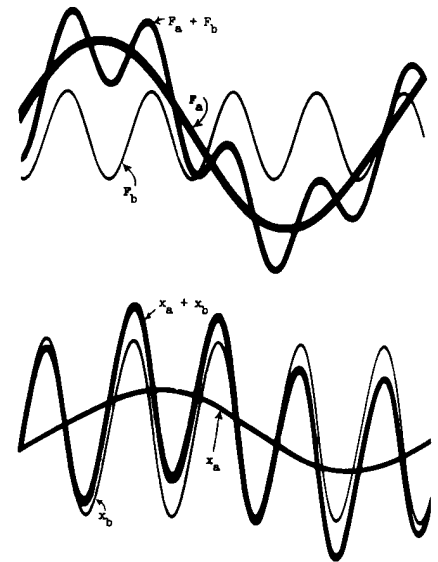


Fig. 25-1. An example of the principle of superposition for linear systems.

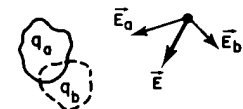


Fig. 25-2. The principle of superposition in electrostatics.

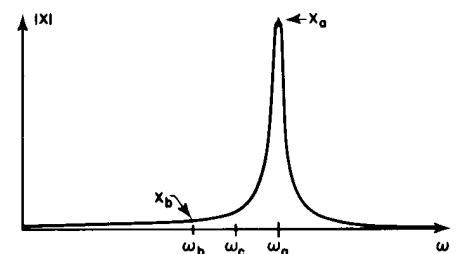


Fig. 25-3. A sharply tuned resonance curve.

the dial, so that the natural frequency of the circuit is shifted, say, to ω_c . In those circumstances we hear neither one station nor the other; we get silence, provided there is no other station at frequency ω_c . If we keep on changing the capacitance until the resonance curve is at ω_b , then of course we hear the other station. That is how radio tuning works; it is again the principle of superposition, combined with a resonant response.*

To conclude this discussion, let us describe qualitatively what happens if we proceed further in analyzing a linear problem with a given force, when the force is quite complicated. Out of the many possible procedures, there are two especially useful general ways that we can solve the problem. One is this: suppose that we can solve it for special known forces, such as sine waves of different frequencies. We know it is child's play to solve it for sine waves. So we have the so-called "child's play" cases. Now the question is whether our very complicated force can be represented as the sum of two or more "child's play" forces. In Fig. 25-1 we already had a fairly complicated curve, and of course we can make it more complicated still if we add in more sine waves. So it is certainly possible to obtain very complicated curves. And, in fact, the reverse is also true: practically every curve can be obtained by adding together *infinite numbers* of sine waves of different wavelengths (or frequencies) for each one of which we know the answer. We just have to know how much of each sine wave to put in to make the given F , and then our answer, x , is the corresponding sum of the F sine waves, each multiplied by its effective ratio of x to F . This method of solution is called the method of *Fourier transforms* or *Fourier analysis*. We are not going to actually carry out such an analysis just now; we only wish to describe the idea involved.

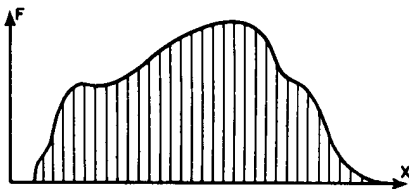


Fig. 25-4. A complicated force may be treated as a succession of sharp impulses.

Another way in which our complicated problem can be solved is the following very interesting one. Suppose that, by some tremendous mental effort, it were possible to solve our problem for a special force, namely an *impulse*. The force is quickly turned on and then off; it is all over. Actually we need only solve for an impulse of some unit strength, any other strength can be gotten by multiplication by an appropriate factor. We know that the response x for an impulse is a damped oscillation. Now what can we say about some other force, for instance a force like that of Fig. 25-4?

Such a force can be likened to a succession of blows with a hammer. First there is no force, and all of a sudden there is a steady force—impulse, impulse, impulse, impulse, . . . and then it stops. In other words, we imagine the continuous force to be a series of impulses, very close together. Now, we know the result for an impulse, so the result for a whole series of impulses will be a whole series of damped oscillations: it will be the curve for the first impulse, and then (slightly later) we add to that the curve for the second impulse, and the curve for the third impulse, and so on. Thus we can represent, mathematically, the complete solution for arbitrary functions if we know the answer for an impulse. We get the answer for any other force simply by integrating. This method is called the *Green's function method*. A Green's function is a response to an impulse, and the method of analyzing any force by putting together the response of impulses is called the *Green's function method*.

The physical principles involved in both of these schemes are so simple, involving just the linear equation, that they can be readily understood, but the *mathematical* problems that are involved, the complicated integrations and so on, are a little too advanced for us to attack right now. You will most likely return to this some day when you have had more practice in mathematics. But the *idea* is very simple indeed.

Finally, we make some remarks on why *linear* systems are so important. The answer is simple: because we can solve them! So most of the time we solve linear

* In modern superheterodyne receivers the actual operation is more complex. The amplifiers are all tuned to a fixed frequency (called IF frequency) and an oscillator of variable tunable frequency is combined with the input signal in a *nonlinear* circuit to produce a new frequency (the difference of signal and oscillator frequency) equal to the IF frequency, which is then amplified. This will be discussed in Chapter 50.

problems. Second (and most important), it turns out that *the fundamental laws of physics are often linear*. The Maxwell equations for the laws of electricity are linear, for example. The great laws of quantum mechanics turn out, so far as we know, to be linear equations. *That* is why we spend so much time on linear equations: because if we understand linear equations, we are ready, in principle, to understand a lot of things.

We mention another situation where linear equations are found. When displacements are small, many functions can be *approximated* linearly. For example, if we have a simple pendulum, the correct equation for its motion is

$$d^2\theta/dt^2 = -(g/L) \sin \theta. \quad (25.9)$$

This equation can be solved by elliptic functions, but the easiest way to solve it is numerically, as was shown in Chapter 9 on Newton's Laws of Motion. A non-linear equation cannot be solved, ordinarily, any other way *but* numerically. Now for small θ , $\sin \theta$ is practically equal to θ , and we have a linear equation. It turns out that there are many circumstances where small effects are linear: for the example here the swing of a pendulum through small arcs. As another example, if we pull a little bit on a spring, the force is proportional to the extension. If we pull hard, we break the spring, and the force is a completely different function of the distance! Linear equations are important. In fact they are *so* important that perhaps fifty percent of the time we are solving linear equations in physics and in engineering.

25-3 Oscillations in linear systems

Let us now review the things we have been talking about in the past few chapters. It is very easy for the physics of oscillators to become obscured by the mathematics. The physics is actually very simple, and if we may forget the mathematics for a moment we shall see that we can understand almost everything that happens in an oscillating system. First, if we have only the spring and the weight, it is easy to understand why the system oscillates—it is a consequence of inertia. We pull the mass down and the force pulls it back up; as it passes zero, which is the place it likes to be, it cannot just suddenly stop; because of its momentum it keeps on going and swings to the other side, and back and forth. So, if there were no friction, we would surely *expect* an oscillatory motion, and indeed we get one. But if there is even a little bit of friction, then on the return cycle, the swing will not be quite as high as it was the first time.

Now what happens, cycle by cycle? That depends on the kind and amount of friction. Suppose that we could concoct a kind of friction force that always remains in the same proportion to the other forces, of inertia and in the spring, as the amplitude of oscillation varies. In other words, for smaller oscillations the friction should be weaker than for big oscillations. Ordinary friction does not have this property, so a special kind of friction must be carefully invented for the very purpose of creating a friction that is directly proportional to the velocity—so that for big oscillations it is stronger and for small oscillations it is weaker. If we happen to have that kind of friction, then at the end of each successive cycle the system is in the same condition as it was at the start, except a little bit smaller. All the forces are smaller in the same proportion: the spring force is reduced, the inertial effects are lower because the accelerations are now weaker, and the friction is less too, by our careful design. When we actually have that kind of friction, we find that each oscillation is exactly the same as the first one, except reduced in amplitude. If the first cycle dropped the amplitude, say, to 90 percent of what it was at the start, the next will drop it to 90 percent of 90 percent, and so on: *the sizes of the oscillations are reduced by the same fraction of themselves in every cycle*. An exponential function is a curve which does just that. It changes by the same factor in each equal interval of time. That is to say, if the amplitude of one cycle, relative to the preceding one, is called a , then the amplitude of the next is a^2 , and of the next, a^3 . So the amplitude is some constant raised to a power equal to the number of cycles traversed:

$$A = A_0 a^n. \quad (25.10)$$

But of course $n \sim t$, so it is perfectly clear that the general solution will be some kind of an oscillation, sine or cosine ωt , times an amplitude which goes as b^t more or less. But b can be written as e^{-c} , if b is positive and less than 1. So this is why the solution looks like $e^{-ct} \cos \omega t$. It is very simple.

What happens if the friction is not so artificial; for example, ordinary rubbing on a table, so that the friction force is a certain constant amount, and is independent of the size of the oscillation that reverses its direction each half-cycle? Then the equation is no longer linear, it becomes hard to solve, and must be solved by the numerical method given in Chapter 2, or by considering each half-cycle separately. The numerical method is the most powerful method of all, and can solve any equation. It is only when we have a simple problem that we can use mathematical analysis.

Mathematical analysis is not the grand thing it is said to be; it solves only the simplest possible equations. As soon as the equations get a little more complicated, just a shade—they cannot be solved analytically. But the numerical method, which was advertised at the beginning of the course, can take care of any equation of physical interest.

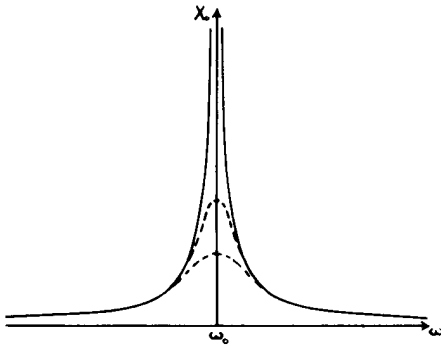


Fig. 25-5. Resonance curves with various amounts of friction present.

Next, what about the resonance curve? Why is there a resonance? First, imagine for a moment that there is no friction, and we have something which could oscillate by itself. If we tapped the pendulum just right each time it went by, of course we could make it go like mad. But if we close our eyes and do not watch it, and tap at arbitrary equal intervals, what is going to happen? Sometimes we will find ourselves tapping when it is going the wrong way. When we happen to have the timing just right, of course, each tap is given at just the right time, and so it goes higher and higher and higher. So without friction we get a curve which looks like the solid curve in Fig. 25-5 for different frequencies. Qualitatively, we understand the resonance curve; in order to get the exact shape of the curve it is probably just as well to do the mathematics. The curve goes toward infinity as $\omega \rightarrow \omega_0$, where ω_0 is the natural frequency of the oscillator.

Now suppose there is a little bit of friction; then when the displacement of the oscillator is small, the friction does not affect it much; the resonance curve is the same, except when we are near resonance. Instead of becoming infinite near resonance, the curve is only going to get so high that the work done by our tapping each time is enough to compensate for the loss of energy by friction during the cycle. So the top of the curve is rounded off—it does not go to infinity. If there is more friction, the top of the curve is rounded off still more. Now someone might say, “I thought the widths of the curves depended on the friction.” That is because the curve is usually plotted so that the top of the curve is called one unit. However, the mathematical expression is even simpler to understand if we just plot all the curves on the same scale; then all that happens is that the friction cuts down the top! If there is less friction, we can go farther up into that little pinnacle before the friction cuts it off, so it looks relatively narrow. That is, the higher the peak of the curve, the narrower the width at half the maximum height.

Finally, we take the case where there is an enormous amount of friction. It turns out that if there is too much friction, the system does not oscillate at all. The energy in the spring is barely able to move it against the frictional force, and so it slowly oozes down to the equilibrium point.

25-4 Analogs in physics

The next aspect of this review is to note that masses and springs are not the only linear systems; there are others. In particular, there are electrical systems called linear circuits, in which we find a complete analog to mechanical systems. We did not learn exactly *why* each of the objects in an electrical circuit works in the way it does—that is not to be understood at the present moment; we may assert it as an experimentally verifiable fact that they behave as stated.

For example, let us take the simplest possible circumstance. We have a piece of wire, which is just a resistance, and we have applied to it a difference in potential, V . Now the V means this: if we carry a charge q through the wire from one terminal

to another terminal, the work done is qV . The higher the voltage difference, the more work was done when the charge, as we say, “falls” from the high potential end of the terminal to the low potential end. So charges release energy in going from one end to the other. Now the charges do not simply fly from one end straight to the other end; the atoms in the wire offer some resistance to the current, and this resistance obeys the following law for almost all ordinary substances: if there is a current I , that is, so and so many charges per second tumbling down, the number per second that comes tumbling through the wire is proportional to how hard we push them—in other words, proportional to how much voltage there is:

$$V = IR = R(dq/dt). \quad (25.11)$$

The coefficient R is called the *resistance*, and the equation is called Ohm’s Law. The unit of resistance is the ohm; it is equal to one volt per ampere. In mechanical situations, to get such a frictional force in proportion to the velocity is difficult; in an electrical system it is very easy, and this law is extremely accurate for most metals.

We are often interested in how much work is done per second, the power loss, or the energy liberated by the charges as they tumble down the wire. When we carry a charge q through a voltage V , the work is qV , so the work done per second would be $V(dq/dt)$, which is the same as VI , or also $IR \cdot I = I^2R$. This is called the *heating loss*—this is how much heat is generated in the resistance per second, by the conservation of energy. It is this heat that makes an ordinary incandescent light bulb work.

Of course, there are other interesting properties of mechanical systems, such as the mass (inertia), and it turns out that there is an electrical analog to inertia also. It is possible to make something called an *inductor*, having a property called *inductance*, such that a current, once started through the inductance, *does not want to stop*. It requires a voltage in order to change the current! If the current is constant, there is no voltage across an inductance. DC circuits do not know anything about inductance; it is only when we *change* the current that the effects of inductance show up. The equation is

$$V = L(dI/dt) = L(d^2q/dt^2), \quad (25.12)$$

and the unit of inductance, called the *henry*, is such that one volt applied to an inductance of one henry produces a change of one ampere per second in the current. Equation (25.12) is the analog of Newton’s law for electricity, if we wish: V corresponds to F , L corresponds to m , and I corresponds to velocity! All of the consequent equations for the two kinds of systems will have the same derivations because, in all the equations, we can change any letter to its corresponding analog letter and we get the same equation; everything we deduce will have a correspondence in the two systems.

Now what electrical thing corresponds to the mechanical spring, in which there was a force proportional to the stretch? If we start with $F = kx$ and replace $F \rightarrow V$ and $x \rightarrow q$, we get $V = \alpha q$. It turns out that there *is* such a thing, in fact it is the only one of the three circuit elements we can really understand, because we did study a pair of parallel plates, and we found that if there were a charge of certain equal, opposite amounts on each plate, the electric field between them would be proportional to the size of the charge. So the work done in moving a unit charge across the gap from one plate to the other is precisely proportional to the charge. This work is the *definition* of the voltage difference, and it is the line integral of the electric field from one plate to another. It turns out, for historical reasons, that the constant of proportionality is not called C , but $1/C$. It could have been called C , but it was not. So we have

$$V = q/C. \quad (25.13)$$

The unit of capacitance, C , is the farad; a charge of one coulomb on each plate of a one-farad capacitor yields a voltage difference of one volt.

There are our analogies, and the equation corresponding to the oscillating circuit becomes the following, by direct substitution of L for m , q for x , etc:

$$m(d^2x/dt^2) + \gamma m(dx/dt) + kx = F, \quad (25.14)$$

$$L(d^2q/dt^2) + R(dq/dt) + q/C = V. \quad (25.15)$$

Now everything we learned about (25.14) can be transformed to apply to (25.15). Every *consequence* is the same; so much the same that there is a brilliant thing we can do.

Suppose we have a mechanical system which is quite complicated, not just one mass on a spring, but several masses on several springs, all hooked together. What do we do? Solve it? Perhaps; but look, we can make an *electrical* circuit which will have the same equations as the thing we are trying to analyze! For instance, if we wanted to analyze a mass on a spring, why can we not build an electrical circuit in which we use an inductance proportional to the mass, a resistance proportional to the corresponding $m\gamma$, $1/C$ proportional to k , all in the same ratio? Then, of course, this electrical circuit will be the exact analog of our mechanical one, in the sense that whatever q does, in response to V (V also is made to correspond to the forces that are acting), so the x would do in response to the force! So if we have a complicated thing with a whole lot of interconnecting elements, we can interconnect a whole lot of resistances, inductances, and capacitances, to *imitate* the mechanically complicated system. What is the advantage to that? One problem is just as hard (or as easy) as the other, because they are exactly equivalent. The advantage is not that it is any easier to solve the *mathematical equations* after we discover that we have an electrical circuit (although that *is* the method used by electrical engineers!), but instead, the real reason for looking at the analog is that it is easier to *make* the electrical circuit, and to *change* something in the system.

Suppose we have designed an automobile, and want to know how much it is going to shake when it goes over a certain kind of bumpy road. We build an electrical circuit with inductances to represent the inertia of the wheels, spring constants as capacitances to represent the springs of the wheels, and resistors to represent the shock absorbers, and so on for the other parts of the automobile. Then we need a bumpy road. All right, we apply a *voltage* from a generator, which represents such and such a kind of bump, and then look at how the left wheel jiggles by measuring the charge on some capacitor. Having measured it (it is easy to do), we find that it is bumping too much. Do we need more shock absorber, or less shock absorber? With a complicated thing like an automobile, do we actually change the shock absorber, and solve it all over again? No!, we simply turn a dial; dial number ten is shock absorber number three, so we put in more shock absorber. The bumps are worse—all right, we try less. The bumps are still worse; we change the stiffness of the spring (dial 17), and we adjust all these things *electrically*, with merely the turn of a knob.

This is called an *analog computer*. It is a device which imitates the problem that we want to solve by making another problem, which has the same equation, but in another circumstance of nature, and which is easier to build, to measure, to adjust, and to destroy!

25-5 Series and parallel impedances

Finally, there is an important item which is not quite in the nature of review. This has to do with an electrical circuit in which there is more than one circuit element. For example, when we have an inductor, a resistor, and a capacitor connected as in Fig. 24-2, we note that all the charge went through every one of the three, so that the current in such a singly connected thing is the same at all points along the wire. Since the current is the same in each one, the voltage across R is IR , the voltage across L is $L(dI/dt)$, and so on. So, the total voltage drop is the sum of these, and this leads to Eq. (25.15). Using complex numbers, we found that we could solve the equation for the steady-state motion in response to a

sinusoidal force. We thus found that $\hat{V} = \hat{Z}\hat{I}$. Now \hat{Z} is called the *impedance* of this particular circuit. It tells us that if we apply a sinusoidal voltage, \hat{V} , we get a current \hat{I} .

Now suppose we have a more complicated circuit which has two pieces, which by themselves have certain impedances, \hat{Z}_1 and \hat{Z}_2 , and we put them in *series* (Fig. 25-6a) and apply a voltage. What happens? It is now a little more complicated, but if \hat{I} is the current through \hat{Z}_1 , the voltage difference across \hat{Z}_1 is $\hat{V}_1 = \hat{I}\hat{Z}_1$; similarly, the voltage across \hat{Z}_2 is $\hat{V}_2 = \hat{I}\hat{Z}_2$. *The same current goes through both.* Therefore the total voltage is the sum of the voltages across the two sections and is equal to $\hat{V} = \hat{V}_1 + \hat{V}_2 = (\hat{Z}_1 + \hat{Z}_2)\hat{I}$. This means that the voltage on the complete circuit can be written $\hat{V} = \hat{I}\hat{Z}_s$ where the \hat{Z}_s of the combined system in series is the sum of the two \hat{Z} 's of the separate pieces:

$$\hat{Z}_s = \hat{Z}_1 + \hat{Z}_2. \quad (25.16)$$

This is not the only way things may be connected. We may also connect them in another way, called a *parallel* connection (Fig. 25-6b). Now we see that a given voltage across the terminals, if the connecting wires are perfect conductors, is effectively applied to both of the impedances, and will cause currents in each independently. Therefore the current through \hat{Z}_1 is equal to $\hat{I}_1 = \hat{V}/\hat{Z}_1$. The current in \hat{Z}_2 is $\hat{I}_2 = \hat{V}/\hat{Z}_2$. It is the *same voltage*. Now the total current which is supplied to the terminals is the *sum* of the currents in the two sections: $\hat{I} = \hat{V}/\hat{Z}_1 + \hat{V}/\hat{Z}_2$. This can be written as

$$\hat{V} = \frac{\hat{I}}{(\hat{I}/\hat{Z}_1) + (\hat{I}/\hat{Z}_2)} = \hat{I}\hat{Z}_p.$$

Thus

$$1/\hat{Z}_p = 1/\hat{Z}_1 + 1/\hat{Z}_2. \quad (25.17)$$

More complicated circuits can sometimes be simplified by taking pieces of them, working out the succession of impedances of the pieces, and combining the circuit together step by step, using the above rules. If we have any kind of circuit with many impedances connected in all kinds of ways, and if we include the voltages in the form of little generators having no impedance (when we pass charge through it, the generator adds a voltage V), then the following principles apply: (1) At any junction, the sum of the currents into a junction is zero. That is, all the current which comes in must come back out. (2) If we carry a charge around any loop, and back to where it started, the net work done is zero. These rules are called *Kirchhoff's laws* for electrical circuits. Their systematic application to complicated circuits often simplifies the analysis of such circuits. We mention them here in conjunction with Eqs. (25.16) and (25.17), in case you have already come across such circuits that you need to analyze in laboratory work. They will be discussed again in more detail next year.

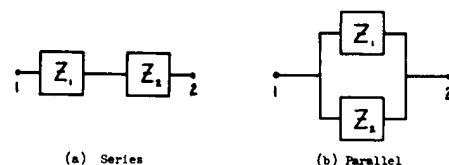


Fig. 25-6. Two impedances, connected in series and in parallel.

Optics: The Principle of Least Time

26-1 Light

This is the first of a number of chapters on the subject of *electromagnetic radiation*. Light, with which we see, is only one small part of a vast spectrum of the same kind of thing, the various parts of this spectrum being distinguished by different values of a certain quantity which varies. This variable quantity could be called the “wavelength.” As it varies in the visible spectrum, the light apparently changes color from red to violet. If we explore the spectrum systematically, from long wavelengths toward shorter ones, we would begin with what are usually called *radiowaves*. Radiowaves are technically available in a wide range of wavelengths, some even longer than those used in regular broadcasts; regular broadcasts have wavelengths corresponding to about 500 meters. Then there are the so-called “short waves,” i.e., radar waves, millimeter waves, and so on. There are no actual boundaries between one range of wavelengths and another, because nature did not present us with sharp edges. The number associated with a given name for the waves are only approximate and, of course, so are the names we give to the different ranges.

Then, a long way down through the millimeter waves, we come to what we call the *infrared*, and thence to the visible spectrum. Then going in the other direction, we get into a region which is called the *ultraviolet*. Where the ultraviolet stops, the x-rays begin, but we cannot define precisely where this is; it is roughly at 10^{-8} m, or 10^{-2} μ . These are “soft” x-rays; then there are ordinary x-rays and very hard x-rays; then γ -rays, and so on, for smaller and smaller values of this dimension called the wavelength.

Within this vast range of wavelengths, there are three or more regions of approximation which are especially interesting. In one of these, a condition exists in which the wavelengths involved are very small compared with the dimensions of the equipment available for their study; furthermore, the photon energies, using the quantum theory, are small compared with the energy sensitivity of the equipment. Under these conditions we can make a rough first approximation by a method called *geometrical optics*. If, on the other hand, the wavelengths are comparable to the dimensions of the equipment, which is difficult to arrange with visible light but easier with radiowaves, and if the photon energies are still negligibly small, then a very useful approximation can be made by studying the behavior of the waves, still disregarding the quantum mechanics. This method is based on the *classical theory of electromagnetic radiation*, which will be discussed in a later chapter. Next, if we go to very short wavelengths, where we can disregard the wave character but the photons have a very *large* energy compared with the sensitivity of our equipment, things get simple again. This is the simple *photon* picture, which we will describe only very roughly. The complete picture, which unifies the whole thing into one model, will not be available to us for a long time.

In this chapter our discussion is limited to the geometrical optics region, in which we forget about the wavelength and the photon character of the light, which will all be explained in due time. We do not even bother to say what the light *is*, but just find out *how it behaves* on a large scale compared with the dimensions of interest. All this must be said in order to emphasize the fact that what we are going to talk about is only a very crude approximation; this is one of the chapters that we shall have to “unlearn” again. But we shall very quickly unlearn it, because we shall almost immediately go on to a more accurate method.

26-1 Light

26-2 Reflection and refraction

26-3 Fermat's principle of least time

26-4 Applications of Fermat's principle

26-5 A more precise statement of Fermat's principle

26-6 How it works

Although geometrical optics is just an approximation, it is of very great importance technically and of great interest historically. We shall present this subject more historically than some of the others in order to give some idea of the development of a physical theory or physical idea.

First, light is, of course, familiar to everybody, and has been familiar since time immemorial. Now one problem is, by what process do we *see* light? There have been many theories, but it finally settled down to one, which is that there is something which enters the eye—which bounces off objects into the eye. We have heard that idea so long that we accept it, and it is almost impossible for us to realize that very intelligent men have proposed contrary theories—that something comes out of the eye and feels for the object, for example. Some other important observations are that, as light goes from one place to another, it goes in *straight lines*, if there is nothing in the way, and that the rays do not seem to interfere with one another. That is, light is crisscrossing in all directions in the room, but the light that is passing across our line of vision does not affect the light that comes to us from some object. This was once a most powerful argument against the corpuscular theory; it was used by Huygens. If light were like a lot of arrows shooting along, how could other arrows go through them so easily? Such philosophical arguments are not of much weight. One could always say that light is made up of arrows which go through each other!

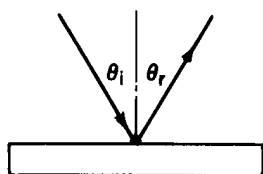


Fig. 26-1. The angle of incidence is equal to the angle of reflection.

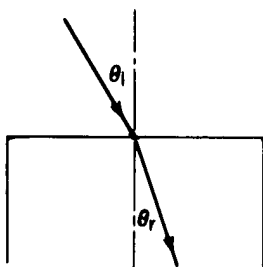


Fig. 26-2. A light ray is refracted when it passes from one medium into another.

26-2 Reflection and refraction

The discussion above gives enough of the basic *idea* of geometrical optics—now we have to go a little further into the quantitative features. Thus far we have light going only in straight lines between two points; now let us study the behavior of light when it hits various materials. The simplest object is a mirror, and the law for a mirror is that when the light hits the mirror, it does not continue in a straight line, but bounces off the mirror into a new straight line, which changes when we change the inclination of the mirror. The question for the ancients was, what is the relation between the two angles involved? This is a very simple relation, discovered long ago. The light striking a mirror travels in such a way that the two angles, between each beam and the mirror, are equal. For some reason it is customary to measure the angles from the normal to the mirror surface. Thus the so-called law of reflection is

$$\theta_i = \theta_r. \quad (26.1)$$

That is a simple enough proposition, but a more difficult problem is encountered when light goes from one medium into another, for example from air into water; here also, we see that it does not go in a straight line. In the water the ray is at an inclination to its path in the air; if we change the angle θ_i so that it comes down more nearly vertically, then the angle of “breakage” is not as great. But if we tilt the beam of light at quite an angle, then the deviation angle is very large. The question is, what is the relation of one angle to the other? This also puzzled the ancients for a long time, and here they never found the answer! It is, however, one of the few places in all of Greek physics that one may find any experimental results listed. Claudius Ptolemy made a list of the angle in water for each of a number of different angles in air. Table 26-1 shows the angles in the air, in degrees, and the corresponding angle as measured in the water. (Ordinarily it is said that Greek scientists never did any experiments. But it would be impossible to obtain this table of values without knowing the right law, except by experiment. It should be noted, however, that these do not represent independent careful measurements for each angle but only some numbers interpolated from a few measurements, for they all fit perfectly on a parabola.)

This, then, is one of the important steps in the development of physical law: first we observe an effect, then we measure it and list it in a table; then we try to find the *rule* by which one thing can be connected with another. The above numerical table was made in 140 A.D., but it was not until 1621 that someone finally found the rule connecting the two angles! The rule, found by Willebrord

Table 26-1

Angle in air	Angle in water
10°	8°
20°	15-1/2°
30°	22-1/2°
40°	29°
50°	35°
60°	40-1/2°
70°	45-1/2°
80°	50°

Snell, a Dutch mathematician, is as follows: if θ_i is the angle in air and θ_r is the angle in the water, then it turns out that the sine of θ_i is equal to some constant multiple of the sine of θ_r :

$$\sin \theta_i = n \sin \theta_r. \quad (26.2)$$

For water the number n is approximately 1.33. Equation (26.2) is called *Snell's law*; it permits us to *predict* how the light is going to bend when it goes from air into water. Table 26-2 shows the angles in air and in water according to Snell's law. Note the remarkable agreement with Ptolemy's list.

Table 26-2

Angle in air	Angle in water
10°	7-1/2°
20°	15°
30°	22°
40°	29°
50°	35°
60°	40-1/2°
70°	45°
80°	48°

26-3 Fermat's principle of least time

Now in the further development of science, we want more than just a formula. First we have an observation, then we have numbers that we measure, then we have a law which summarizes all the numbers. But the real *glory* of science is that we can find a way of thinking such that the law is *evident*.

The first way of thinking that made the law about the behavior of light evident was discovered by Fermat in about 1650, and it is called *the principle of least time*, or *Fermat's principle*. His idea is this: that out of all possible paths that it might take to get from one point to another, light takes the path which requires the *shortest time*.

Let us first show that this is true for the case of the mirror, that this simple principle contains both the law of straight-line propagation and the law for the mirror. So, we are growing in our understanding! Let us try to find the solution to the following problem. In Fig. 26-3 are shown two points, A and B , and a plane mirror, MM' . What is the way to get from A to B in the shortest time? The answer is to go straight from A to B ! But if we add the extra rule that the light has to *strike the mirror* and come back in the shortest time, the answer is not so easy. One way would be to go as quickly as possible to the mirror and then go to B , on the path ADB . Of course, we then have a long path DB . If we move over a little to the right, to E , we slightly increase the first distance, but we greatly *decrease* the second one, and so the total path length, and therefore the travel time, is less. How can we find the point C for which the time is the shortest? We can find it very nicely by a geometrical trick.

We construct on the other side of MM' an artificial point B' , which is the same distance below the plane MM' as the point B is above the plane. Then we draw the line EB' . Now because BFM is a right angle and $BF = FB'$, EB is equal to EB' . Therefore the sum of the two distances, $AE + EB$, which is proportional to the time it will take if the light travels with constant velocity, is also the sum of the two lengths $AE + EB'$. Therefore the problem becomes, when is the sum of these two lengths the least? The answer is easy: when the line goes through point C as a *straight line* from A to B' ! In other words, we have to find the point where we go toward the artificial point, and that will be the correct one. Now if ACB' is a straight line, then angle BCF is equal to angle $B'CF$ and thence to angle ACM . Thus the statement that the angle of incidence equals the angle of reflection is equivalent to the statement that the light goes to the mirror in such a way that it comes back to the point B' in the *least possible time*. Originally, the statement was made by Hero of Alexandria that the light travels in such a way that it goes to the mirror and to the other point in the shortest possible *distance*, so it is not a modern theory. It was this that inspired Fermat to suggest to himself that perhaps refraction operated on a similar basis. But for refraction, light obviously does not use the path of shortest *distance*, so Fermat tried the idea that it takes the shortest *time*.

Before we go on to analyze refraction, we should make one more remark about the mirror. If we have a source of light at the point B and it sends light toward the mirror, then we see that the light which goes to A from the point B comes to A in exactly the same manner as it would have come to A if there were an object at B' , and *no* mirror. Now of course the eye detects only the light which enters it physically, so if we have an object at B and a mirror which makes the light come

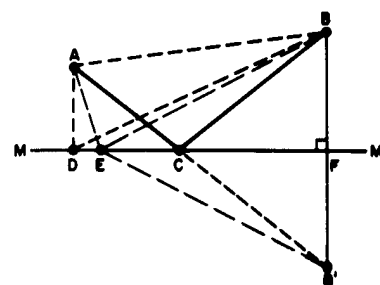


Fig. 26-3. Illustration of the principle of least time.

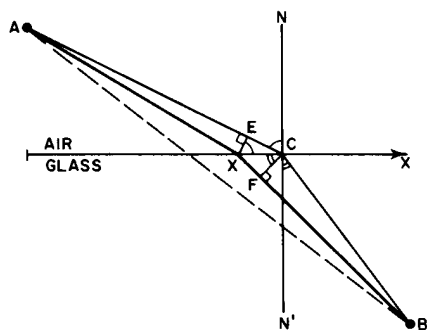


Fig. 26-4. Illustration of Fermat's principle for refraction.

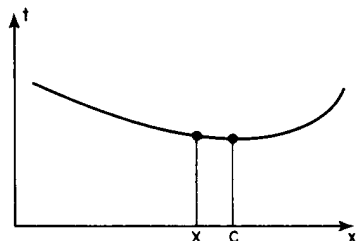


Fig. 26-5. The minimum time corresponds to point C, but nearby points correspond to nearly the same time.

into the eye in exactly the same manner as it would have come into the eye if the object were at B' , then the eye-brain system interprets that, assuming it does not know too much, as *being* an object at B' . So the illusion that there is an object behind the mirror is merely due to the fact that the light which is entering the eye is entering in exactly the same manner, physically, as it would have entered had there *been* an object back there (except for the dirt on the mirror, and our knowledge of the existence of the mirror, and so on, which is corrected in the brain).

Now let us demonstrate that the principle of least time will give Snell's law of refraction. We must, however, make an assumption about the speed of light in water. We shall assume that the speed of light in water is lower than the speed of light in air by a certain factor, n .

In Fig. 26-4, our problem is again to go from A to B in the *shortest time*. To illustrate that the best thing to do is not just to go in a straight line, let us imagine that a beautiful girl has fallen out of a boat, and she is screaming for help in the water at point B . The line marked X is the shoreline. We are at point A on land, and we see the accident, and we can run and can also swim. But we can run faster than we can swim. What do we do? Do we go in a straight line? (Yes, no doubt!) However, by using a little more intelligence we would realize that it would be advantageous to travel a little greater distance on land in order to decrease the distance in the water, because we go so much slower in the water. (Following this line of reasoning out, we would say the right thing to do is to *compute* very carefully what should be done!) At any rate, let us try to show that the final solution to the problem is the path ACB , and that this path takes the shortest time of all possible ones. If it is the shortest path, that means that if we take any other, it will be longer. So, if we were to plot the time it takes against the position of point X , we would get a curve something like that shown in Fig. 26-5, where point C corresponds to the shortest of all possible times. This means that if we move the point X to points *near* C , in the first approximation there is essentially *no change* in time because the slope is zero at the bottom of the curve. So our way of finding the law will be to consider that we move the place by a very small amount, and to demand that there be essentially no change in time. (Of course there is an infinitesimal change of a *second* order; we ought to have a positive increase for displacements in either direction from C .) So we consider a nearby point X and we calculate how long it would take to go from A to B by the two paths, and compare the new path with the old path. It is very easy to do. We want the difference, of course, to be nearly zero if the distance XC is short. First, look at the path on land. If we draw a perpendicular XE , we see that this path is shortened by the amount EC . Let us say we gain by not having to go that extra distance. On the other hand, in the water, by drawing a corresponding perpendicular, CF , we find that we have to go the extra distance XF , and that is what we lose. Or, in *time*, we gain the time it would have taken to go the distance EC , but we lose the time it would have taken to go the distance XF . Those times must be equal since, in the first approximation, there is to be no change in time. But supposing that in the water the speed is $1/n$ times as fast as in air, then we must have

$$EC = n \cdot XF. \quad (26.3)$$

Therefore we see that when we have the right point, $XC \sin EXC = n \cdot XC \sin XCF$ or, cancelling the common hypotenuse length XC and noting that

$$EXC = ECN = \theta_i \quad \text{and} \quad XCF = BCN' = \theta_r,$$

we have

$$\sin \theta_i = n \sin \theta_r. \quad (26.4)$$

So we see that to get from one point to another in the least time when the ratio of speeds is n , the light should enter at such an angle that the ratio of the sines of the angles θ_i and θ_r is the ratio of the speeds in the two media.

26-4 Applications of Fermat's principle

Now let us consider some of the interesting consequences of the principle of least time. First is the principle of reciprocity. If to go from A to B we have found the path of the least time, then to go in the opposite direction (assuming that light goes at the same speed in any direction), the shortest time will be the same path, and therefore, if light can be sent one way, it can be sent the other way.

An example of interest is a glass block with plane parallel faces, set at an angle to a light beam. Light, in going through the block from a point A to a point B (Fig. 26-6) does not go through in a straight line, but instead it decreases the time in the block by making the angle in the block less inclined, although it loses a little bit in the air. The beam is simply displaced parallel to itself because the angles in and out are the same.

A third interesting phenomenon is the fact that when we see the sun setting, it is already below the horizon! It does not *look* as though it is below the horizon, but it is (Fig. 26-7). The earth's atmosphere is thin at the top and dense at the bottom. Light travels more slowly in air than it does in a vacuum, and so the light of the sun can get to point S beyond the horizon more quickly if, instead of just going in a straight line, it avoids the dense regions where it goes slowly by getting through them at a steeper tilt. When it appears to go below the horizon, it is actually already well below the horizon. Another example of this phenomenon is the mirage that one often sees while driving on hot roads. One sees "water" on the road, but when he gets there, it is as dry as the desert! The phenomenon is the following. What we are really seeing is the sky light "reflected" on the road: light from the sky, heading for the road, can end up in the eye, as shown in Fig. 26-8. Why? The air is very hot just above the road but it is cooler up higher. Hotter air is more expanded than cooler air and is thinner, and this decreases the speed of light less. That is to say, light goes faster in the hot region than in the cool region. Therefore, instead of the light deciding to come in the straightforward way, it also has a least-time path by which it goes into the region where it goes faster for awhile, in order to save time. So, it can go in a curve.

As another important example of the principle of least time, suppose that we would like to arrange a situation where we have all the light that comes out of one point, P , collected back together at another point, P' (Fig. 26-9). That means, of course, that the light can go in a straight line from P to P' . That is all right. But how can we arrange that not only does it go straight, but also so that the light starting out from P toward Q also ends up at P' ? We want to bring all the light back to what we call a *focus*. How? If the light always takes the path of least time, then certainly it should not want to go over all these other paths. The only way that the light can be perfectly satisfied to take several adjacent paths is to make those times *exactly equal*! Otherwise, it would select the one of least time. Therefore the problem of making a focusing system is merely to arrange a device so that it takes the same time for the light to go on *all* the different paths!

This is easy to do. Suppose that we had a piece of glass in which light goes slower than it does in the air (Fig. 26-10). Now consider a ray which goes in air in the path PQP' . That is a longer path than from P directly to P' and no doubt takes a longer time. But if we were to insert a piece of glass of just the right thickness (we shall later figure out how thick) it might exactly compensate the excess time that it would take the light to go at an angle! In those circumstances we can arrange that the time the light takes to go straight through is the same as the time it takes to go in the path PQP' . Likewise, if we take a ray $PRR'P'$ which is partly inclined, it is not quite as long as PQP' , and we do not have to compensate as much as for the straight one, but we do have to compensate somewhat. We end up with a piece of glass that looks like Fig. 26-10. With this shape, all the light which comes from P will go to P' . This, of course, is well known to us, and we call such a device a converging *lens*. In the next chapter we shall actually calculate what shape the lens has to have to make a perfect focus.

Take another example: suppose we wish to arrange some mirrors so that the light from P always goes to P' (Fig. 26-11). On any path, it goes to some mirror

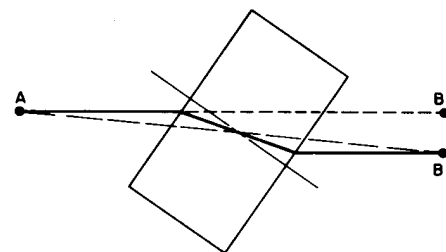


Fig. 26-6. A beam of light is offset as it passes through a transparent block.

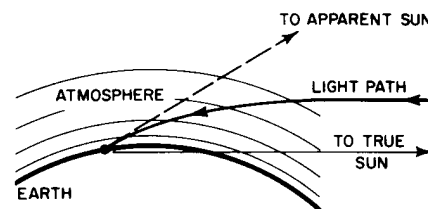


Fig. 26-7. Near the horizon, the apparent sun is higher than the true sun by about $1/2$ degree.

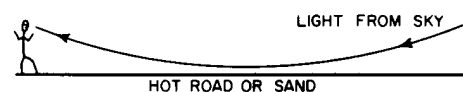


Fig. 26-8. A mirage.

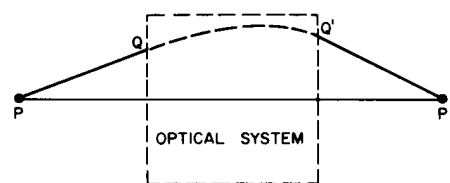


Fig. 26-9. An optical "black box."

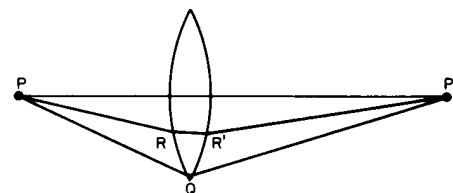


Fig. 26-10. A focusing optical system.

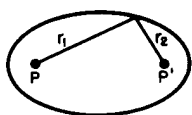


Fig. 26-11. An ellipsoidal mirror.

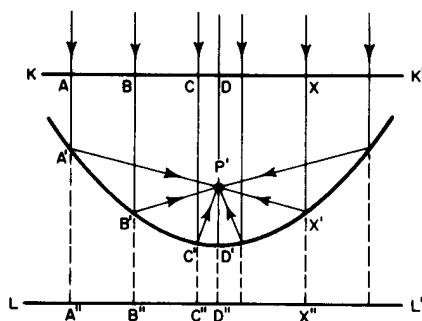


Fig. 26-12. A paraboloidal mirror.

and comes back, and all times must be equal. Here the light always travels in air, so the time and the distance are proportional. Therefore the statement that all the times are the same is the same as the statement that the total distance is the same. Thus the sum of the two distances r_1 and r_2 must be a constant. An *ellipse* is that curve which has the property that the sum of the distances from two points is a constant for every point on the ellipse; thus we can be sure that the light from one focus will come to the other.

The same principle works for gathering the light of a star. The great 200-inch Palomar telescope is built on the following principle. Imagine a star billions of miles away; we would like to cause all the light that comes in to come to a focus. Of course we cannot draw the rays that go all the way up to the star, but we still want to check whether the times are equal. Of course we know that when the various rays have arrived at some plane KK' , perpendicular to the rays, all the times in this plane are equal (Fig. 26-12). The rays must then come down to the mirror and proceed toward P' in equal times. That is, we must find a curve which has the property that the sum of the distances $XX' + X'P'$ is a constant, no matter where X is chosen. An easy way to find it is to extend the length of the line XX' down to a plane LL' . Now if we arrange our curve so that $A'A'' = A'P'$, $B'B'' = B'P'$, $C'C'' = C'P'$, and so on, we will have our curve, because then of course, $AA' + A'P' = AA' + A'A''$ will be constant. Thus our curve is the locus of all points equidistant from a line and a point. Such a curve is called a *parabola*; the mirror is made in the shape of a parabola.

The above examples illustrate the principle upon which such optical devices can be designed. The exact curves can be calculated using the principle that, to focus perfectly, the travel times must be exactly equal for all light rays, as well as being less than for any other nearby path.

We shall discuss these focusing optical devices further in the next chapter; let us now discuss the further development of the theory. When a new theoretical principle is developed, such as the principle of least time, our first inclination might be to say, "Well, that is very pretty; it is delightful; but the question is, does it help at all in understanding the physics?" Someone may say, "Yes, look at how many things we can now understand!" Another says, "Very well, but I can understand mirrors, too. I need a curve such that every tangent plane makes equal angles with the two rays. I can figure out a lens, too, because every ray that comes to it is bent through an angle given by Snell's law." Evidently the statement of least time and the statement that angles are equal on reflection, and that the sines of the angles are proportional on refraction, are the same. So is it merely a philosophical question, or one of beauty? There can be arguments on both sides.

However, the importance of a powerful principle is that *it predicts new things*.

It is easy to show that there are a number of new things predicted by Fermat's principle. First, suppose that there are *three* media, glass, water, and air, and we perform a refraction experiment and measure the index n for one medium against another. Let us call n_{12} the index of air (1) against water (2); n_{13} the index of air (1) against glass (3). If we measured water against glass, we should find another index, which we shall call n_{23} . But there is no *a priori* reason why there should be any connection between n_{12} , n_{13} , and n_{23} . On the other hand, according to the idea of least time, there *is* a definite relationship. The index n_{12} is the ratio of two things, the speed in air to the speed in water; n_{13} is the ratio of the speed in air to the speed in glass; n_{23} is the ratio of the speed in water to the speed in glass. Therefore we cancel out the air, and get

$$n_{23} = \frac{v_2}{v_3} = \frac{v_1/v_3}{v_1/v_2} = \frac{n_{13}}{n_{12}}. \quad (26.5)$$

In other words, we *predict* that the index for a new pair of materials can be obtained from the indexes of the individual materials, both against air or against vacuum. So if we measure the speed of light in all materials, and from this get a single number for each material, namely its index relative to vacuum, called n_i

(n_1 is the speed in air relative to the speed in vacuum, etc.), then our formula is easy. The index for any two materials i and j is

$$n_{ij} = \frac{v_i}{v_j} = \frac{n_j}{n_i}. \quad (26.6)$$

Using only Snell's law, there is no basis for a prediction of this kind.* But of course this prediction works. The relation (26.5) was known very early, and was a very strong argument for the principle of least time.

Another argument for the principle of least time, another prediction, is that if we *measure* the speed of light in water, it will be lower than in air. This is a prediction of a completely different type. It is a brilliant prediction, because all we have so far measured are *angles*; here we have a theoretical prediction which is quite different from the observations from which Fermat deduced the idea of least time. It turns out, in fact, that the speed in water *is* slower than the speed in air, by just the proportion that is needed to get the right index!

26-5 A more precise statement of Fermat's principle

Actually, we must make the statement of the principle of least time a little more accurately. It was not stated correctly above. It is *incorrectly* called the principle of least time and we have gone along with the incorrect description for convenience, but we must now see what the correct statement is. Suppose we had a mirror as in Fig. 26-3. What makes the light think it has to go to the mirror? The path of *least* time is clearly AB . So some people might say, "Sometimes it is a maximum time." It is *not* a maximum time, because certainly a curved path would take a still longer time! The correct statement is the following: a ray going in a certain particular path has the property that if we make a small change (say a one percent shift) in the ray in any manner whatever, say in the location at which it comes to the mirror, or the shape of the curve, or anything, there will be *no* first-order change in the time; there will be only a *second*-order change in the time. In other words, the principle is that light takes a path such that there are many other paths nearby which take almost exactly the *same* time.

The following is another difficulty with the principle of least time, and one which people who do not like this kind of a theory could never stomach. With Snell's theory we can "understand" light. Light goes along, it sees a surface, it bends because it does something at the surface. The idea of causality, that it goes from one point to another, and another, and so on, is easy to understand. But the principle of least time is a completely different philosophical principle about the way nature works. Instead of saying it is a causal thing, that when we do one thing, something else happens, and so on, it says this: we set up the situation, and *light* decides which is the shortest time, or the extreme one, and chooses that path. But *what* does it do, *how* does it find out? Does it *smell* the nearby paths, and check them against each other? The answer is, yes, it does, in a way. That is the feature which is, of course, not known in geometrical optics, and which is involved in the idea of *wavelength*; the wavelength tells us approximately how far away the light must "smell" the path in order to check it. It is hard to demonstrate this fact on a large scale with light, because the wavelengths are so terribly short. But with radiowaves, say 3-cm waves, the distances over which the radiowaves are checking are larger. If we have a source of radiowaves, a detector, and a slit, as in Fig. 26-13, the rays of course go from S to D because it is a straight line, and if we close down the slit it is all right—they still go. But now if we move the detector aside to D' , the waves will not go through the wide slit from S to D' , because they check several paths nearby, and say, "No, my friend, those all correspond to different times." On the other hand, if we *prevent* the radiation from checking the paths by closing the slit down to a very narrow crack, then there is but one path available, and the

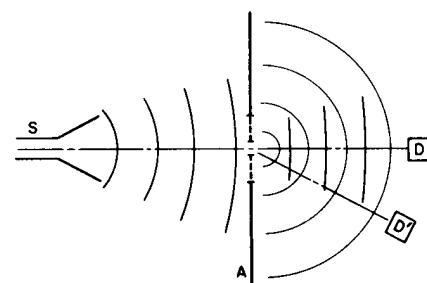


Fig. 26-13. The passage of radio-waves through a narrow slit.

* Although it can be deduced if the additional assumption is made that adding a layer of one substance to the surface of another does not change the eventual angle of refraction in the latter material.

radiation takes it! With a narrow slit, more radiation reaches D' than reaches it with a wide slit!

One can do the same thing with light, but it is hard to demonstrate on a large scale. The effect can be seen under the following simple conditions. Find a small, bright light, say an unfrosted bulb in a street light far away or the reflection of the sun in a curved automobile bumper. Then put two fingers in front of one eye, so as to look through the crack, and squeeze the light to zero very gently. You will see that the image of the light, which was a little dot before, becomes quite elongated, and even stretches into a long line. The reason is that the fingers are very close together, and the light which is supposed to come in a straight line is spread out at an angle, so that when it comes into the eye it comes in from several directions. Also you will notice, if you are very careful, side maxima, a lot of fringes along the edges too. Furthermore, the whole thing is colored. All of this will be explained in due time, but for the present it is a demonstration that light does not always go in straight lines, and it is one that is very easily performed.

26-6 How it works

Finally, we give a very crude view of what actually happens, how the whole thing really works, from what we now believe is the correct, quantum-dynamically accurate viewpoint, but of course only qualitatively described. In following the light from A to B in Fig. 26-3, we find that the light does not seem to be in the form of waves at all. Instead the rays seem to be made up of photons, and they actually produce clicks in a photon counter, if we are using one. The brightness of the light is proportional to the average number of photons that come in per second, and what we calculate is the *chance* that a photon gets from A to B , say by hitting the mirror. The *law* for that chance is the following very strange one. Take any path and find the time for that path; then make a complex number, or draw a little complex vector, $pe^{i\theta}$, whose angle θ is proportional to the time. The number of turns per second is the frequency of the light. Now take another path; it has, for instance, a different time, so the vector for it is turned through a different angle—the angle being always proportional to the time. Take *all* the available paths and add on a little vector for each one; then the answer is that the chance of arrival of the photon is proportional to the square of the length of the final vector, from the beginning to the end!

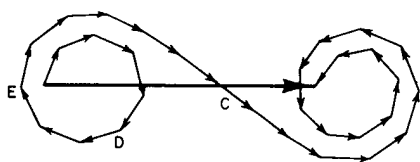


Fig. 26-14. The summation of probability amplitudes for many neighboring paths.

Now let us show how this implies the principle of least time for a mirror. We consider all rays, all possible paths ADB , AEB , ACB , etc., in Fig. 26-3. The path ADB makes a certain small contribution, but the next path, AEB , takes a quite different time, so its angle θ is quite different. Let us say that point C corresponds to minimum time, where if we change the paths the times do not change. So for awhile the times do change, and then they begin to change less and less as we get near point C (Fig. 26-14). So the arrows which we have to add are coming almost exactly at the same angle for awhile near C , and then gradually the time begins to increase again, and the phases go around the other way, and so on. Eventually, we have quite a tight knot. The total probability is the distance from one end to the other, squared. *Almost all of that accumulated probability occurs in the region where all the arrows are in the same direction* (or in the same phase). All the contributions from the paths which have very *different* times as we change the path, cancel themselves out by pointing in different directions. That is why, if we hide the extreme parts of the mirror, it still reflects almost exactly the same, because all we did was to take out a piece of the diagram inside the spiral ends, and that makes only a very small change in the light. So this is the relationship between the ultimate picture of photons with a probability of arrival depending on an accumulation of arrows, and the principle of least time.

Geometrical Optics

27-1 Introduction

In this chapter we shall discuss some elementary applications of the ideas of the previous chapter to a number of practical devices, using the approximation called *geometrical optics*. This is a most useful approximation in the practical design of many optical systems and instruments. Geometrical optics is either very simple or else it is very complicated. By that we mean that we can either study it only superficially, so that we can design instruments roughly, using rules that are so simple that we hardly need deal with them here at all, since they are practically of high school level, or else, if we want to know about the small errors in lenses and similar details, the subject gets so complicated that it is too advanced to discuss here! If one has an actual, detailed problem in lens design, including analysis of aberrations, then he is advised to read about the subject or else simply to trace the rays through the various surfaces (which is what the book tells how to do), using the law of refraction from one side to the other, and to find out where they come out and see if they form a satisfactory image. People have said that this is too tedious, but today, with computing machines, it is the right way to do it. One can set up the problem and make the calculation for one ray after another very easily. So the subject is really ultimately quite simple, and involves no new principles. Furthermore, it turns out that the rules of either elementary or advanced optics are seldom characteristic of other fields, so that there is no special reason to follow the subject very far, with one important exception.

The most advanced and abstract theory of geometrical optics was worked out by Hamilton, and it turns out that this has very important applications in mechanics. It is actually even more important in mechanics than it is in optics, and so we leave Hamilton's theory for the subject of advanced analytical mechanics, which is studied in the senior year or in graduate school. So, appreciating that geometrical optics contributes very little, except for its own sake, we now go on to discuss the elementary properties of simple optical systems on the basis of the principles outlined in the last chapter.

In order to go on, we must have one geometrical formula, which is the following: if we have a triangle with a small altitude h and a long base d , then the diagonal s (we are going to need it to find the difference in time between two different routes) is longer than the base (Fig. 27-1). How much longer? The difference $\Delta = s - d$ can be found in a number of ways. One way is this. We see that $s^2 - d^2 = h^2$, or $(s - d)(s + d) = h^2$. But $s - d = \Delta$, and $s + d \sim 2s$. Thus

$$\Delta \sim h^2/2s. \quad (27.1)$$

This is all the geometry we need to discuss the formation of images by curved surfaces!

27-2 The focal length of a spherical surface

The first and simplest situation to discuss is a single refracting surface, separating two media with different indices of refraction (Fig. 27-2). We leave the case of arbitrary indices of refraction to the student, because *ideas* are always the most important thing, not the specific situation, and the problem is easy enough to do in any case. So we shall suppose that, on the left, the speed is 1 and on the right it is $1/n$, where n is the index of refraction. The light travels more slowly in the glass by a factor n .

27-1 Introduction

27-2 The focal length of a spherical surface

27-3 The focal length of a lens

27-4 Magnification

27-5 Compound lenses

27-6 Aberrations

27-7 Resolving power

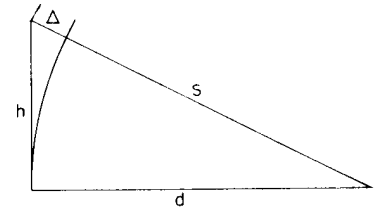


Figure 27-1

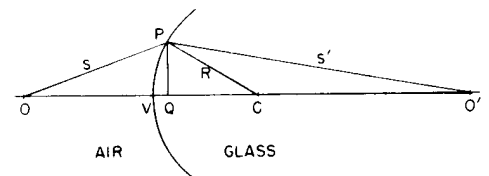


Fig. 27-2. Focusing by a single refracting surface.

Now suppose that we have a point at O , at a distance s from the front surface of the glass, and another point O' at a distance s' inside the glass, and we desire to arrange the curved surface in such a manner that every ray from O which hits the surface, at any point P , will be bent so as to proceed toward the point O' . For that to be true, we have to shape the surface in such a way that the time it takes for the light to go from O to P , that is, the distance OP divided by the speed of light (the speed here is unity), plus $n \cdot O'P$, which is the time it takes to go from P to O' , is equal to a constant independent of the point P . This condition supplies us with an equation for determining the surface. The answer is that the surface is a very complicated fourth-degree curve, and the student may entertain himself by trying to calculate it by analytic geometry. It is simpler to try a special case that corresponds to $s \rightarrow \infty$, because then the curve is a second-degree curve and is more recognizable. It is interesting to compare this curve with the parabolic curve we found for a focusing mirror when the light is coming from infinity.

So the proper surface cannot easily be made—to focus the light from one point to another requires a rather complicated surface. It turns out in practice that we do not try to make such complicated surfaces ordinarily, but instead we make a compromise. Instead of trying to get *all* the rays to come to a focus, we arrange it so that only the rays fairly close to the axis OO' come to a focus. The farther ones may deviate if they want to, unfortunately, because the ideal surface is complicated, and we use instead a spherical surface with the right curvature at the axis. It is so much easier to fabricate a sphere than other surfaces that it is profitable for us to find out what happens to rays striking a spherical surface, supposing that only the rays near the axis are going to be focused perfectly. Those rays which are near the axis are sometimes called *paraxial rays*, and what we are analyzing are the conditions for the focusing of paraxial rays. We shall discuss later the errors that are introduced by the fact that all rays are not always close to the axis.

Thus, supposing P is close to the axis, we drop a perpendicular PQ such that the height PQ is h . For a moment, we imagine that the surface is a plane passing through P . In that case, the time needed to go from O to P would exceed the time from O to Q , and also, the time from P to O' would exceed the time from Q to O' . But that is why the glass must be curved, because the total excess time must be compensated by the delay in passing from V to Q ! Now the *excess* time along route OP is $h^2/2s$, and the excess time on the other route is $nh^2/2s'$. This excess time, which must be matched by the delay in going along VQ , differs from what it would have been in a vacuum, because there is a medium present. In other words, the time to go from V to Q is not as if it were straight in the air, but it is slower by the factor n , so that the excess delay in this distance is then $(n - 1)VQ$. And now, how large is VQ ? If the point C is the center of the sphere and if its radius is R , we see by the same formula that the distance VQ is equal to $h^2/2R$. Therefore we discover that the law that connects the distances s and s' , and that gives us the radius of curvature R of the surface that we need, is

$$(h^2/2s) + (nh^2/2s') = (n - 1)h^2/2R \quad (27.2)$$

or

$$(1/s) + (n/s') = (n - 1)/R. \quad (27.3)$$

If we have a position O and another position O' , and want to focus light from O to O' , then we can calculate the required radius of curvature R of the surface by this formula.

Now it turns out, interestingly, that the same lens, with the same curvature R , will focus for other distances, namely, for any pair of distances such that the sum of the two reciprocals, one multiplied by n , is a constant. Thus a given lens will (so long as we limit ourselves to paraxial rays) focus not only from O to O' , but between an infinite number of other pairs of points, so long as those pairs of points bear the relationship that $1/s + n/s'$ is a constant, characteristic of the lens.

In particular, an interesting case is that in which $s \rightarrow \infty$. We can see from the formula that as one s increases, the other decreases. In other words, if point O

goes out, point O' comes in, and vice versa. As point O goes toward infinity, point O' keeps moving in until it reaches a certain distance, called the *focal length* f' , inside the material. If parallel rays come in, they will meet the axis at a distance f' . Likewise, we could imagine it the other way. (Remember the reciprocity rule: if light will go from O to O' , of course it will also go from O' to O .) Therefore, if we had a light source inside the glass, we might want to know where the focus is. In particular, if the light in the glass were at infinity (same problem) where would it come to a focus outside? This distance is called f . Of course, we can also put it the other way. If we had a light source at f and the light went through the surface, then it would go out as a parallel beam. We can easily find out what f and f' are:

$$n/f' = (n - 1)/R \quad \text{or} \quad f' = Rn/(n - 1), \quad (27.4)$$

$$1/f = (n - 1)/R \quad \text{or} \quad f = R/(n - 1). \quad (27.5)$$

We see an interesting thing: if we divide each focal length by the corresponding index of refraction we get the same result! This theorem, in fact, is general. It is true of any system of lenses, no matter how complicated, so it is interesting to remember. We did not prove here that it is general—we merely noted it for a single surface, but it happens to be true in general that the two focal lengths of a system are related in this way. Sometimes Eq. (27.3) is written in the form

$$1/s + n/s' = 1/f. \quad (27.6)$$

This is more useful than (27.3) because we can measure f more easily than we can measure the curvature and index of refraction of the lens: if we are not interested in designing a lens or in knowing how it got that way, but simply lift it off a shelf, the interesting quantity is f , not the n and the 1 and the R !

Now an interesting situation occurs if s becomes less than f . What happens then? If $s < f$, then $(1/s) > (1/f)$, and therefore s' is negative; our equation says that the light will focus only with a negative value of s' , whatever that means! It does mean something very interesting and very definite. It is still a useful formula, in other words, even when the numbers are negative. What it means is shown in Fig. 27-3. If we draw the rays which are diverging from O , they will be bent, it is true, at the surface, and they will not come to a focus, because O is so close in that they are “beyond parallel.” However, they diverge as if they had come from a point O' outside the glass. This is an apparent image, sometimes called a *virtual image*. The image O' in Fig. 27-2 is called a *real image*. If the light really comes to a point, it is a real image. But if the light *appears* to be coming from a point, a fictitious point different from the original point, it is a virtual image. So when s' comes out negative, it means that O' is on the other side of the surface, and everything is all right.

Now consider the interesting case where R is equal to infinity; then we have $(1/s) + (n/s') = 0$. In other words, $s' = -ns$, which means that if we look from a dense medium into a rare medium and see a point in the rare medium, it appears to be deeper by a factor n . Likewise, we can use the same equation backwards, so that if we look into a plane surface at an object that is at a certain distance inside the dense medium, it will appear as though the light is coming from not as far back (Fig. 27-4). When we look at the bottom of a swimming pool from above, it does not look as deep as it really is, by a factor $3/4$, which is the reciprocal of the index of refraction of water.

We could go on, of course, to discuss the spherical mirror. But if one appreciates the ideas involved, he should be able to work it out for himself. Therefore we leave it to the student to work out the formula for the spherical mirror, but we mention that it is well to adopt certain conventions concerning the distances involved:

- (1) The object distance s is positive if the point O is to the left of the surface.
- (2) The image distance s' is positive if the point O' is to the right of the surface.
- (3) The radius of curvature of the surface is positive if the center is to the right of the surface.

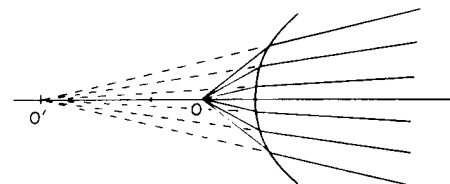


Fig. 27-3. A virtual image.

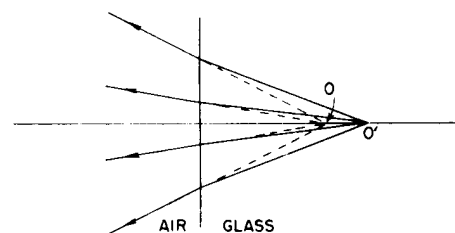


Fig. 27-4. A plane surface re-images the light from O' to O .

In Fig. 27-2, for example, s , s' , and R are all positive; in Fig. 27-3, s and R are positive, but s' is negative. If we had used a concave surface, our formula (27.3) would still give the correct result if we merely make R a negative quantity.

In working out the corresponding formula for a mirror, using the above conventions, you will find that if you put $n = -1$ throughout the formula (27.3) (as though the material behind the mirror had an index -1), the right formula for a mirror results!

Although the derivation of formula (27.3) is simple and elegant, using least time, one can of course work out the same formula using Snell's law, remembering that the angles are so small that the sines of angles can be replaced by the angles themselves.

27-3 The focal length of a lens

Now we go on to consider another situation, a very practical one. Most of the lenses that we use have two surfaces, not just one. How does this affect matters? Suppose that we have two surfaces of different curvature, with glass filling the space between them (Fig. 27-5). We want to study the problem of focusing from a point O to an alternate point O' . How can we do that? The answer is this: First, use formula (27.3) for the first surface, forgetting about the second surface. This will tell us that the light which was diverging from O will appear to be converging or diverging, depending on the sign, from some other point, say O'' . Now we consider a new problem. We have a different surface, between glass and air, in which rays are converging toward a certain point O' . Where will they actually converge? We use the same formula again! We find that they converge at O''' . Thus, if necessary, we can go through 75 surfaces by just using the same formula in succession, from one to the next!

There are some rather high-class formulas that would save us considerable energy in the few times in our lives that we might have to chase the light through five surfaces, but it is easier just to chase it through five surfaces when the problem arises than it is to memorize a lot of formulas, because it may be we will never have to chase it through any surfaces at all!

In any case, the principle is that when we go through one surface we find a new position, a new focal point, and then take that point as the starting point for the next surface, and so on. In order to actually do this, since on the second surface we are going from n to 1 rather than from 1 to n , and since in many systems there is more than one kind of glass, so that there are indices n_1, n_2, \dots , we really need a generalization of formula (27.3) for a case where there are two different indices, n_1 and n_2 , rather than only n . Then it is not difficult to prove that the general form of (27.3) is

$$(n_1/s) + (n_2/s') = (n_2 - n_1)/R. \quad (27.7)$$

Particularly simple is the special case in which the two surfaces are very close together—so close that we may ignore small errors due to the thickness. If we draw the lens as shown in Fig. 27-6, we may ask this question: How must the lens be built so as to focus light from O to O' ? Suppose the light comes exactly to the edge of the lens, at point P . Then the excess time in going from O to O' is $(n_1 h^2/2s) + (n_1 h^2/2s')$, ignoring for a moment the presence of the thickness T of glass of index n_2 . Now, to make the time for the direct path equal to that for the path OPO' , we have to use a piece of glass whose thickness T at the center is such that the delay introduced in going through this thickness is enough to compensate for the excess time above. Therefore the thickness of the lens at the center must be given by the relationship

$$(n_1 h^2/2s) + (n_1 h^2/2s') = (n_2 - n_1)T. \quad (27.8)$$

We can also express T in terms of the radii R_1 and R_2 of the two surfaces. Paying attention to our convention (3), we thus find, for $R_1 < R_2$ (a convex lens),

$$T = (h^2/2R_1) - (h^2/2R_2). \quad (27.9)$$

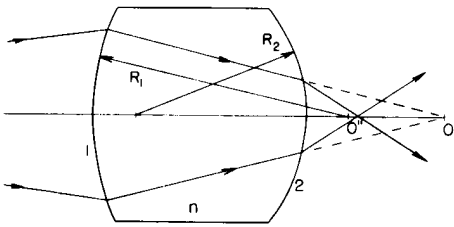


Fig. 27-5. Image formation by a two-surface lens.

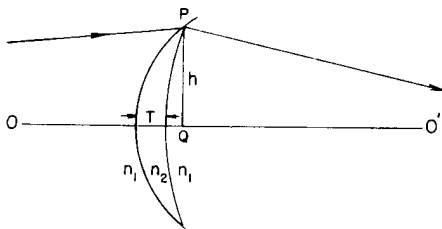


Fig. 27-6. A thin lens with two positive radii.

Therefore, we finally get

$$(n_1/s) + (n_1/s') = (n_2 - n_1)(1/R_1 - 1/R_2). \quad (27.10)$$

Now we note again that if one of the points is at infinity, the other will be at a point which we will call the focal length f . The focal length f is given by

$$1/f = (n - 1)(1/R_1 - 1/R_2), \quad (27.11)$$

where $n = n_2/n_1$.

Now, if we take the opposite case, where s goes to infinity, we see that s' is at the focal length f' . This time the focal lengths are equal. (This is another special case of the general rule that the ratio of the two focal lengths is the ratio of the indices of refraction in the two media in which the rays focus. In this particular optical system, the initial and final indices are the same, so the two focal lengths are equal.)

Forgetting for a moment about the actual formula for the focal length, if we bought a lens that somebody designed with certain radii of curvature and a certain index, we could measure the focal length, say, by seeing where a point at infinity focuses. Once we had the focal length, it would be better to write our equation in terms of the focal length directly, and the formula then is

$$(1/s) + (1/s') = 1/f. \quad (27.12)$$

Now let us see how the formula works and what it implies in different circumstances. First, it implies that if s or s' is infinite the other one is f . That means that parallel light focuses at a distance f , and this in effect *defines* f . Another interesting thing it says is that both points move in the same direction. If one moves to the right, the other does also. Another thing it says is that s and s' are equal if they are both equal to $2f$. In other words, if we want a symmetrical situation, we find that they will both focus at a distance $2f$.

27-4 Magnification

So far we have discussed the focusing action only for points on the axis. Now let us discuss also the imaging of objects not exactly on the axis, but a little bit off, so that we can understand the properties of *magnification*. When we set up a lens so as to focus light from a small filament onto a "point" on a screen, we notice that on the screen we get a "picture" of the same filament, except of a larger or smaller size than the true filament. This must mean that the light comes to a focus from *each point* of the filament. In order to understand this a little better, let us analyze the thin lens system shown schematically in Fig. 27-7. We know the following facts:

- (1) Any ray that comes in parallel on one side proceeds toward a certain particular point called the focus on the other side, at a distance f from the lens.
- (2) Any ray that arrives at the lens from the focus on one side comes out parallel to the axis on the other side.

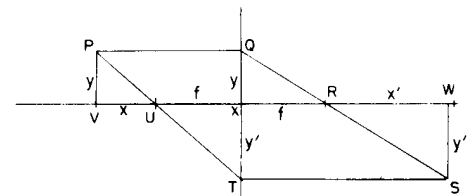


Fig. 27-7. The geometry of imaging by a thin lens.

This is all we need to establish formula (27.12) by geometry, as follows: Suppose we have an object at some distance x from the focus; let the height of the object be y . Then we know that one of the rays, namely PQ , will be bent so as to pass through the focus R on the other side. Now if the lens will focus point P at all, we can find out where if we find out where just one other ray goes, because the new focus will be where the two intersect again. We need only use our ingenuity to find the exact direction of *one* other ray. But we remember that a parallel ray goes through the focus and *vice versa*: a ray which goes through the focus will come out parallel! So we draw ray PT through U . (It is true that the actual rays which are doing the focusing may be much more limited than the two we have drawn, but they are harder to figure, so we make believe that we can make this ray.) Since it would come out parallel, we draw TS parallel to XW . The intersection S is the point we need. This will determine the correct place and the correct height. Let

us call the height y' and the distance from the focus, x' . Now we may derive a lens formula. Using the similar triangles PVU and TXU , we find

$$\frac{y'}{f} = \frac{y}{x}. \quad (27.13)$$

Similarly, from triangles SWR and QXR , we get

$$\frac{y'}{x'} = \frac{y}{f}. \quad (27.14)$$

Solving each for y'/y , we find that

$$\frac{y'}{y} = \frac{x'}{f} = \frac{f}{x}. \quad (27.15)$$

Equation (27.15) is the famous lens formula; in it is everything we need to know about lenses: It tells us the magnification, y'/y , in terms of the distances and the focal lengths. It also connects the two distances x and x' with f :

$$xx' = f^2, \quad (27.16)$$

which is a much neater form to work with than Eq. (27.12). We leave it to the student to demonstrate that if we call $s = x + f$ and $s' = x' + f$, Eq. (27.12) is the same as Eq. (27.16).

27-5 Compound lenses

Without actually deriving it, we shall briefly describe the general result when we have a number of lenses. If we have a system of several lenses, how can we possibly analyze it? That is easy. We start with some object and calculate where its image is for the first lens, using formula (27.16) or (27.12) or any other equivalent formula, or by drawing diagrams. So we find an image. Then we treat this image as the source for the next lens, and use the second lens with whatever its focal length is to again find an image. We simply chase the thing through the succession of lenses. That is all there is to it. It involves nothing new in principle, so we shall not go into it. However, there is a very interesting net result of the effects of any sequence of lenses on light that starts and ends up in the same medium, say air. Any optical instrument—a telescope or a microscope with any number of lenses and mirrors—has the following property: There exist two planes, called the *principal planes* of the system (these planes are often fairly close to the first surface of the first lens and the last surface of the last lens), which have the following properties: (1) If light comes into the system parallel from the first side, it comes out at a certain focus, at a distance from the *second* principal plane equal to the focal length, just as though the system were a thin lens situated at this plane. (2) If parallel light comes in the other way, it comes to a focus at the same distance f from the *first* principal plane, again as if a thin lens were situated there. (See Fig. 27-8.)

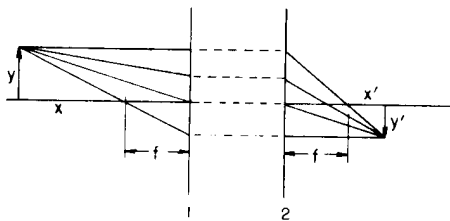


Fig. 27-8. Illustration of the principal planes of an optical system.

Of course, if we measure the distances x and x' , and y and y' as before, the formula (27.16) that we have written for the thin lens is absolutely general, provided that we measure the focal length from the principal planes and not from the center of the lens. It so happens that for a thin lens the principal planes are coincident. It is just as though we could take a thin lens, slice it down the middle, and separate it, and not notice that it was separated. Every ray that comes in pops out immediately on the other side of the second plane from the same point as it went into the first plane! The principal planes and the focal length may be found either by experiment or by calculation, and then the whole set of properties of the optical system are described. It is very interesting that the result is not complicated when we are all finished with such a big, complicated optical system.

27-6 Aberrations

Before we get too excited about how marvelous lenses are, we must hasten to add that there are also serious limitations, because of the fact that we have limited ourselves, strictly speaking, to paraxial rays, the rays near the axis. A real lens having a finite size will, in general, exhibit *aberrations*. For example, a ray that is on the axis, of course, goes through the focus; a ray that is very close to the axis will still come to the focus very well. But as we go farther out, the ray begins to deviate from the focus, perhaps by falling short, and a ray striking near the top edge comes down and misses the focus by quite a wide margin. So, instead of getting a point image, we get a smear. This effect is called *spherical aberration*, because it is a property of the spherical surfaces we use in place of the right shape. This could be remedied, for any specific object distance, by re-forming the shape of the lens surface, or perhaps by using several lenses arranged so that the aberrations of the individual lenses tend to cancel each other.

Lenses have another fault: light of different colors has different speeds, or different indices of refraction, in the glass, and therefore the focal length of a given lens is different for different colors. So if we image a white spot, the image will have colors, because when we focus for the red, the blue is out of focus, or vice versa. This property is called *chromatic aberration*.

There are still other faults. If the object is off the axis, then the focus really isn't perfect any more, when it gets far enough off the axis. The easiest way to verify this is to focus a lens and then tilt it so that the rays are coming in at a large angle from the axis. Then the image that is formed will usually be quite crude, and there may be no place where it focuses well. There are thus several kinds of errors in lenses that the optical designer tries to remedy by using many lenses to compensate each other's errors.

How careful do we have to be to eliminate aberrations? Is it possible to make an absolutely perfect optical system? Suppose we had built an optical system that is supposed to bring light exactly to a point. Now, arguing from the point of view of least time, can we find a condition on how perfect the system has to be? The system will have some kind of an entrance opening for the light. If we take the farthest ray from the axis that can come to the focus (if the system is perfect, of course), the times for all rays are exactly equal. But nothing is perfect, so the question is, how wrong can the time be for this ray and not be worth correcting any further? That depends on how perfect we want to make the image. But suppose we want to make the image as perfect as it possibly can be made. Then, of course, our impression is that we have to arrange that every ray takes as nearly the same time as possible. But it turns out that this is not true, that beyond a certain point we are trying to do something that is too fine, because the theory of geometrical optics does not work!

Remember that the principle of least time is not an accurate formulation, unlike the principle of conservation of energy or the principle of conservation of momentum. The principle of least time is only an *approximation*, and it is interesting to know how much error can be allowed and still not make any apparent difference. The answer is that if we have arranged that between the maximal ray—the worst ray, the ray that is farthest out—and the central ray, the difference in time is less than about the period that corresponds to one oscillation of the light, then there is no use improving it any further. Light is an oscillatory thing with a definite frequency that is related to the wavelength, and if we have arranged that the time difference for different rays is less than about a period, there is no use going any further.

27-7 Resolving power

Another interesting question—a very important technical question with all optical instruments—is how much *resolving power* they have. If we build a microscope, we want to see the objects that we are looking at. That means, for instance, that if we are looking at a bacterium with a spot on each end, we want to *see* that

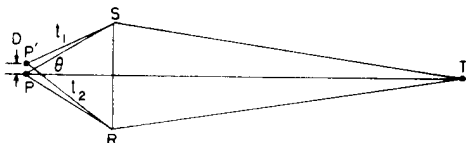


Fig. 27-9. The resolving power of an optical system.

there are two dots when we magnify them. One might think that all we have to do is to get enough magnification—we can always add another lens, and we can always magnify again and again, and with the cleverness of designers, all the spherical aberrations and chromatic aberrations can be cancelled out, and there is no reason why we cannot keep on magnifying the image. So the limitations of a microscope are not that it is impossible to build a lens that magnifies more than 2000 diameters. We can build a system of lenses that magnifies 10,000 diameters, but we *still* could not see two points that are too close together because of the limitations of geometrical optics, because of the fact that least time is not precise.

To discover the rule that determines how far apart two points have to be so that at the image they appear as separate points can be stated in a very beautiful way associated with the time it takes for different rays. Suppose that we disregard the aberrations now, and imagine that for a particular point P (Fig. 27-9) all the rays from object to image T take exactly the same time. (It is not true, because it is not a perfect system, but that is another problem.) Now take another nearby point, P' , and ask whether its image will be distinct from T . In other words, whether we can make out the difference between them. Of course, according to geometrical optics, there should be two point images, but what we see may be rather smeared and we may not be able to make out that there are two points. The condition that the second point is focused in a distinctly different place from the first one is that the two times for the extreme rays $P'ST$ and $P'RT$ on each side of the big opening of the lenses to go from one end to the other, must *not* be equal from the two possible object points to a given image point. Why? Because, if the times were equal, of course both would *focus* at the same point. So the times are not going to be equal. But by how much do they have to differ so that we can say that both do *not* come to a common focus, so that we can distinguish the two image points? The general rule for the resolution of any optical instrument is this: two different point sources can be resolved only if one source is focused at such a point that the times for the maximal rays from the other source to reach that point, as compared with its own true image point, differ by more than one period. It is necessary that the difference in time between the top ray and the bottom ray to the *wrong* focus shall exceed a certain amount, namely, approximately the period of oscillation of the light:

$$t_2 - t_1 > 1/\nu, \quad (27.17)$$

where ν is the frequency of the light (number of oscillations per second; also speed divided by wavelength). If the distance of separation of the two points is called D , and if the opening angle of the lens is called θ , then one can demonstrate that (27.17) is exactly equivalent to the statement that D must exceed $\lambda/n \sin \theta$, where n is the index of refraction at P and λ is the wavelength. The smallest things that we can see are therefore approximately the wavelength of light. A corresponding formula exists for telescopes, which tells us the smallest difference in angle between two stars that can just be distinguished.*

* The angle is about λ/D , where D is the lens diameter. Can you see why?

Electromagnetic Radiation

28-1 Electromagnetism

The most dramatic moments in the development of physics are those in which great syntheses take place, where phenomena which previously had appeared to be different are suddenly discovered to be but different aspects of the same thing. The history of physics is the history of such syntheses, and the basis of the success of physical science is mainly that we are *able* to synthesize.

Perhaps the most dramatic moment in the development of physics during the 19th century occurred to J. C. Maxwell one day in the 1860's, when he combined the laws of electricity and magnetism with the laws of the behavior of light. As a result, the properties of light were partly unravelled—that old and subtle stuff that is so important and mysterious that it was felt necessary to arrange a special creation for it when writing Genesis. Maxwell could say, when he was finished with his discovery, “Let there be electricity and magnetism, and there is light!”

For this culminating moment there was a long preparation in the gradual discovery and unfolding of the laws of electricity and magnetism. This story we shall reserve for detailed study next year. However, the story is, briefly, as follows. The gradually discovered properties of electricity and magnetism, of electric forces of attraction and repulsion, and of magnetic forces, showed that although these forces were rather complex, they all fell off inversely as the square of the distance. We know, for example, that the simple coulomb law for stationary charges is that the electric force field varies inversely as the square of the distance. As a consequence, for sufficiently great distances there is very little influence of one system of charges on another. Maxwell noted that the equations or the laws that had been discovered up to this time were mutually inconsistent when he tried to put them all together, and in order for the whole system to be consistent, he had to add another term to his equations. With this new term there came an amazing prediction, which was that a part of the electric and magnetic fields would fall off much more slowly with the distance than the inverse square, namely, inversely as the first power of the distance! And so he realized that electric currents in one place can affect other charges far away, and he predicted the basic effects with which we are familiar today—radio transmission, radar, and so on.

It seems a miracle that someone talking in Europe can, with mere electrical influences, be heard thousands of miles away in Los Angeles. How is it possible? It is because the fields do not vary as the inverse square, but only inversely as the first power of the distance. Finally, then, even light itself was recognized to be electric and magnetic influences extending over vast distances, generated by an almost incredibly rapid oscillation of the electrons in the atoms. All these phenomena we summarize by the word *radiation* or, more specifically, *electromagnetic radiation*, there being one or two other kinds of radiation also. Almost always, radiation means electromagnetic radiation.

And thus is the universe knit together. The atomic motions of a distant star still have sufficient influence at this great distance to set the electrons in our eye in motion, and so we know about the stars. If this law did not exist, we would all be literally in the dark about the exterior world! And the electric surgings in a galaxy five billion light years away—which is the farthest object we have found so far—can still influence in a significant and detectable way the currents in the great “dish” in front of a radio telescope. And so it is that we see the stars and the galaxies.

28-1 Electromagnetism

28-2 Radiation

28-3 The dipole radiator

28-4 Interference

This remarkable phenomenon is what we shall discuss in the present chapter. At the beginning of this course in physics we outlined a broad picture of the world, but we are now better prepared to understand some aspects of it, and so we shall now go over some parts of it again in greater detail. We begin by describing the position of physics at the end of the 19th century. All that was then known about the fundamental laws can be summarized as follows.

First, there were laws of forces: one force was the law of gravitation, which we have written down several times; the force on an object of mass m , due to another of mass M , is given by

$$\mathbf{F} = GmMe_r/r^2, \quad (28.1)$$

where \mathbf{e}_r is a unit vector directed from m to M , and r is the distance between them.

Next, the laws of electricity and magnetism, as known at the end of the 19th century, are these: the electrical forces acting on a charge q can be described by two fields, called \mathbf{E} and \mathbf{B} , and the velocity \mathbf{v} of the charge q , by the equation

$$\mathbf{F} = q(\mathbf{E} + \mathbf{v} \times \mathbf{B}). \quad (28.2)$$

To complete this law, we have to say what the formulas for \mathbf{E} and \mathbf{B} are in a given circumstance: if a number of charges are present, \mathbf{E} and the \mathbf{B} are each the sum of contributions, one from each individual charge. So if we can find the \mathbf{E} and \mathbf{B} produced by a single charge, we need only to add all the effects from all the charges in the universe to get the total \mathbf{E} and \mathbf{B} ! This is the principle of superposition.

What is the formula for the electric and magnetic field produced by one individual charge? It turns out that this is very complicated, and it takes a great deal of study and sophistication to appreciate it. But that is not the point. We write down the law now only to impress the reader with the beauty of nature, so to speak, i.e., that it is *possible* to summarize all the fundamental knowledge on one page, with notations that he is now familiar with. This law for the fields of an individual charge is complete and accurate, so far as we know (except for quantum mechanics) but it looks rather complicated. We shall not study all the pieces now; we only write it down to give an impression, to show that it can be written, and so that we can see ahead of time roughly what it looks like. As a matter of fact, the most *useful* way to write the correct laws of electricity and magnetism is not the way we shall now write them, but involves what are called *field equations*, which we shall learn about next year. But the mathematical notations for these are different and new, and so we write the law in an inconvenient form for calculation, but in notations that we now know.

The electric field, \mathbf{E} , is given by

$$\mathbf{E} = \frac{-q}{4\pi\epsilon_0} \left[\frac{\mathbf{e}_{r'}}{r'^2} + \frac{r'}{c} \frac{d}{dt} \left(\frac{\mathbf{e}_{r'}}{r'^2} \right) + \frac{1}{c^2} \frac{d^2}{dt^2} \mathbf{e}_{r'} \right]. \quad (28.3)$$

What do the various terms tell us? Take the first term, $\mathbf{E} = -q\mathbf{e}_{r'}/4\pi\epsilon_0 r'^2$. That, of course, is Coulomb's law, which we already know: q is the charge that is producing the field; $\mathbf{e}_{r'}$ is the unit vector in the direction from the point P where \mathbf{E} is measured, r is the distance from P to q . But, Coulomb's law is wrong. The discoveries of the 19th century showed that influences cannot travel faster than a certain fundamental speed c , which we now call the speed of light. It is not correct that the first term is Coulomb's law, not only because it is not possible to know where the charge is *now* and at what distance it is *now*, but also because the only thing that can affect the field at a given place and time is the behavior of the charges in the *past*. How *far* in the past? The time delay, or *retarded time*, so-called, is the time it takes, at speed c , to get from the charge to the field point P . The delay is r'/c .

So to allow for this time delay, we put a little prime on r , meaning how far away it *was* when the information now arriving at P left q . Just for a moment suppose that the charge carried a light, and that the light could only come to P at the speed c . Then when we look at q , we would not see where it is now, of course, but where it *was* at some earlier time. What appears in our formula is the *apparent*

direction $\theta_{r'}$ —the direction it used to be—the so-called *retarded* direction—and at the *retarded* distance r' . That would be easy enough to understand, too, but it is also wrong. The whole thing is much more complicated.

There are several more terms. The next term is as though nature were trying to allow for the fact that the effect is retarded, if we might put it very crudely. It suggests that we should calculate the delayed coulomb field and add a correction to it, which is its rate of change times the time delay that we use. Nature seems to be attempting to guess what the field at the present time is going to be, by taking the rate of change and multiplying by the time that is delayed. But we are not yet through. There is a third term—the second derivative, with respect to t , of the unit vector in the direction of the charge. Now the formula *is* finished, and that is all there is to the electric field from an arbitrarily moving charge.

The magnetic field is given by

$$\mathbf{B} = -\mathbf{e}_{r'} \times \mathbf{E}/c. \quad (28.4)$$

We have written these down only for the purpose of showing the beauty of nature or, in a way, the power of mathematics. We do not pretend to understand *why* it is possible to write so much in such a small space, but (28.3) and (28.4) contain the machinery by which electric generators work, how light operates, all the phenomena of electricity and magnetism. Of course, to complete the story we also need to know something about the behavior of the materials involved—the properties of matter—which are not described properly by (28.3).

To finish with our description of the world of the 19th century we must mention one other great synthesis which occurred in that century, one with which Maxwell had a great deal to do also, and that was the synthesis of the phenomena of heat and mechanics. We shall study that subject soon.

What had to be added in the 20th century was that the dynamical laws of Newton were found to be all wrong, and quantum mechanics had to be introduced to correct them. Newton's laws are approximately valid when the scale of things is sufficiently large. These quantum-mechanical laws, combined with the laws of electricity, have only recently been combined to form a set of laws called *quantum electrodynamics*. In addition, there were discovered a number of new phenomena, of which the first was radioactivity, discovered by Becquerel in 1898—he just sneaked it in under the 19th century. This phenomenon of radioactivity was followed up to produce our knowledge of nuclei and new kinds of forces that are not gravitational and not electrical, but new particles with different interactions, a subject which has still not been unravelled.

For those purists who know more (the professors who happen to be reading this), we should add that when we say that (28.3) is a complete expression of the knowledge of electrodynamics, we are not being entirely accurate. There was a problem that was not quite solved at the end of the 19th century. When we try to calculate the field from all the charges *including the charge itself that we want the field to act on*, we get into trouble trying to find the distance, for example, of a charge from itself, and dividing something by that distance, which is zero. The problem of how to handle the part of this field which is generated by the very charge on which we want the field to act is not yet solved today. So we leave it there; we do not have a complete solution to that puzzle yet, and so we shall avoid the puzzle for as long as we can.

28-2 Radiation

That, then, is a summary of the world picture. Now let us use it to discuss the phenomena called radiation. To discuss these phenomena, we must select from Eq. (28.3) only that piece which varies inversely as the distance and not as the square of the distance. It turns out that when we finally do find that piece, it is so simple in its form that it is legitimate to study optics and electrodynamics in an elementary way by taking it as “the law” of the electric field produced by a moving charge far away. We shall take it temporarily as a given law which we will learn about in detail next year.

Of the terms appearing in (28.3), the first one evidently goes inversely as the square of the distance, and the second is only a correction for delay, so it is easy to show that both of them vary inversely as the square of the distance. All of the effects we are interested in come from the third term, which is not very complicated, after all. What this term says is: look at the charge and note the direction of the unit vector (we can project the end of it onto the surface of a unit sphere). As the charge moves around, the unit vector wiggles, and *the acceleration of that unit vector is what we are looking for*. That is all. Thus

$$\mathbf{E} = \frac{-q}{4\pi\epsilon_0 c^2} \frac{d^2 \mathbf{e}_{r'}}{dt^2} \quad (28.5)$$

is a statement of the laws of radiation, because that is the only important term when we get far enough away that the fields are varying inversely as the distance. (The parts that go as the square have fallen off so much that we are not interested in them.)

Now we can go a little bit further in studying (28.5) to see what it means. Suppose a charge is moving in any manner whatsoever, and we are observing it from a distance. We imagine for a moment that in a sense it is “lit up” (although it is light that we are trying to explain); we imagine it as a little white dot. Then we would see this white dot running around. But we don’t see *exactly* how it is running around right *now*, because of the delay that we have been talking about. What counts is how it was moving *earlier*. The unit vector $\mathbf{e}_{r'}$ is pointed toward the apparent position of the charge. Of course, the end of $\mathbf{e}_{r'}$ goes on a slight curve, so that its acceleration has two components. One is the transverse piece, because the end of it goes up and down, and the other is a radial piece because it stays on a sphere. It is easy to demonstrate that the latter is much smaller and varies as the inverse square of r when r is very great. This is easy to see, for when we imagine that we move a given source farther and farther away, then the wiggles of $\mathbf{e}_{r'}$ look smaller and smaller, inversely as the distance, but the radial component of acceleration is varying much more rapidly than inversely as the distance. So for practical purposes all we have to do is project the motion on a plane at unit distance. Therefore we find the following rule: Imagine that we look at the moving charge and that everything we see is delayed—like a painter trying to paint a scene on a screen at a unit distance. A real painter, of course, does not take into account the fact that light is going at a certain speed, but paints the world as he sees it. We want to see what his picture would look like. So we see a dot, representing the charge, moving about in the picture. The acceleration of that dot is proportional to the electric field. That is all—all we need.

Thus Eq. (28.5) is the complete and correct formula for radiation; even relativity effects are all contained in it. However, we often want to apply it to a still simpler circumstance in which the charges are moving only a small distance at a relatively slow rate. Since they are moving slowly, they do not move an appreciable distance from where they start, so that the delay time is practically constant. Then the law is still simpler, because the delay time is fixed. Thus we imagine that the charge is executing a very tiny motion at an effectively constant distance. The delay at the distance r is r/c . Then our rule becomes the following: If the charged object is moving in a very small motion and it is laterally displaced by the distance $x(t)$, then the angle that the unit vector $\mathbf{e}_{r'}$ is displaced is x/r , and since r is practically constant, the x -component of $d^2 \mathbf{e}_{r'}/dt^2$ is simply the acceleration of x itself at an earlier time, and so finally we get the law we want, which is

$$E_x(t) = \frac{-q}{4\pi\epsilon_0 c^2 r} a_x \left(t - \frac{r}{c} \right). \quad (28.6)$$

Only the component of a_x perpendicular to the line of sight is important. Let us see why that is. Evidently, if the charge is moving in and out straight at us, the unit vector in that direction does not wiggle at all, and it has no acceleration. So it is only the sidewise motion which is important, only the acceleration that we see projected on the screen.

28-3 The dipole radiator

As our fundamental “law” of electromagnetic radiation, we are going to assume that (28.6) is true, i.e., that the electric field produced by an accelerating charge which is moving nonrelativistically at a very large distance r approaches that form. The electric field varies inversely as r and is proportional to the acceleration of the charge, projected onto the “plane of sight,” and this acceleration is not today’s acceleration, but the acceleration that it had at an earlier time, the amount of delay being a time, r/c . In the remainder of this chapter we shall discuss this law so that we can understand it better physically, because we are going to use it to understand all of the phenomena of light and radio propagation, such as reflection, refraction, interference, diffraction, and scattering. It is the central law, and is all we need. All the rest of Eq. (28.3) was written down only to set the stage, so that we could appreciate where (28.6) fits and how it comes about.

We shall discuss (28.3) further next year. In the meantime, we shall accept it as true, but not just on a theoretical basis. We may devise a number of experiments which illustrate the character of the law. In order to do so, we need an accelerating charge. It should be a single charge, but if we can make a great many charges move together, all the same way, we know that the field will be the sum of the effects of each of the individual charges; we just add them together. As an example, consider two pieces of wire connected to a generator, as shown in Fig. 28-1. The idea is that the generator makes a potential difference, or a field, which pulls electrons away from piece A and pushes them into B at one moment, and then, an infinitesimal time later, it reverses the effect and pulls the electrons out of B and pumps them back into A ! So in these two wires charges, let us say, are accelerating upward in wire A and upward in wire B for one moment, and a moment later they are accelerating downward in wire A and downward in wire B . The fact that we need two wires and a generator is merely that this is a way of doing it. The net result is that we merely have a charge accelerating up and down as though A and B were one single wire. A wire that is very short compared with the distance light travels in one oscillation period is called an *electric dipole oscillator*. Thus we have the circumstance that we need to apply our law, which tells us that this charge makes an electric field, and so we need an instrument to detect an electric field, and the instrument we use is the same thing—a pair of wires like A and B ! If an electric field is applied to such a device, it will produce a force which will pull the electrons up on both wires or down on both wires. This signal is detected by means of a rectifier mounted between A and B , and a tiny, fine wire carries the information into an amplifier, where it is amplified so we can hear the audio-frequency tone with which the radiofrequency is modulated. When this probe feels an electric field, there will be a loud noise coming out of the loudspeaker, and when there is no electric field driving it, there will be no noise.

Because the room in which the waves we are measuring has other objects in it, our electric field will shake electrons in these other objects; the electric field makes these other charges go up and down, and in going up and down, these also produce an effect on our probe. Thus for a successful experiment we must hold things fairly close together, so that the influences from the walls and from ourselves—the reflected waves—are relatively small. So the phenomena will not turn out to appear to be precisely and perfectly in accord with Eq. (28.6), but will be close enough that we shall be able to appreciate the law.

Now we turn the generator on and hear the audio signal. We find a strong field when the detector D is parallel to the generator G at point 1 (Fig. 28-2). We find the same amount of field also at any other azimuth angle about the axis of G , because it has no directional effects. On the other hand, when the detector is at 3 the field is zero. That is all right, because our formula said that the field should be the acceleration of the charge *projected perpendicular* to the line of sight. Therefore when we look down on G , the charge is moving toward and away from D , and there is no effect. So that checks the first rule, that there is no effect when the charge is moving directly toward us. Secondly, the formula says that the electric field should be perpendicular to r and in the plane of G and r ; so if we put D at 1 but rotate it 90°, we should get no signal. And this is just what we find, the

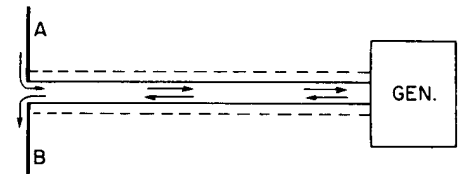


Fig. 28-1. A high-frequency signal generator drives charges up and down on two wires.

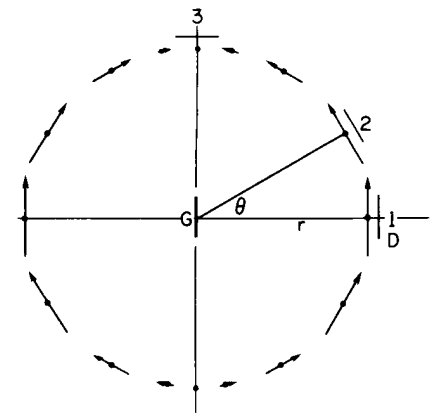


Fig. 28-2. The instantaneous electric field on a sphere centered at a localized, linearly oscillating charge.

electric field is indeed vertical, and not horizontal. When we move D to some intermediate angle, we see that the strongest signal occurs when it is oriented as shown, because although G is vertical, it does not produce a field that is simply parallel to itself—it is the *projection of the acceleration perpendicular to the line of sight* that counts. The signal is weaker at 2 than it is at 1, because of the projection effect.

28-4 Interference

Next, we may test what happens when we have two sources side by side a few centimeters apart (Fig. 28-3). The law is that the two sources should add their effects at point 1 when both of the sources are connected to the same generator and are both moving up and down the same way, so that the total electric field is the sum of the two and is twice as strong as it was before.

Now comes an interesting possibility. Suppose we make the charges in S_1 and S_2 both accelerate up and down, but delay the timing of S_2 so that they are 180° out of phase. Then the field produced by S_1 will be in one direction and the field produced by S_2 will be in the opposite direction at any instant, and therefore we should get *no* effect at point 1. The phase of oscillation is neatly adjustable by means of a pipe which is carrying the signal to S_2 . By changing the length of this pipe we change the time it takes the signal to arrive at S_2 and thus we change the phase of that oscillation. By adjusting this length, we can indeed find a place where there is no more signal left, in spite of the fact that both S_1 and S_2 are moving! The fact that they are both moving can be checked, because if we cut one out, we can see the motion of the other. So the two of them together can produce *zero* if everything is adjusted correctly.

Now, it is very interesting to show that the addition of the two fields is in fact a *vector* addition. We have just checked it for up and down motion, but let us check two nonparallel directions. First, we restore S_1 and S_2 to the same phase; that is, they are again moving together. But now we turn S_1 through 90° , as shown in Fig. 28-4. Now we should have at point 1 the sum of two effects, one of which is vertical and the other horizontal. The electric field is the vector sum of these two in-phase signals—they are both strong at the same time and go through zero together; the total field should be a signal R at 45° . If we turn D to get the maximum noise, it should be at about 45° , and not vertical. And if we turn it at right angles to that direction, we should get zero, which is easy to measure. Indeed, we observe just such behavior!

Now, how about the retardation? How can we demonstrate that the signal is retarded? We could, with a great deal of equipment, measure the time at which it arrives, but there is another, very simple way. Referring again to Fig. 28-3, suppose that S_1 and S_2 are in phase. They are both shaking together, and they produce equal electric fields at point 1. But suppose we go to a certain place 2 which is closer to S_2 and farther from S_1 . Then, in accordance with the principle that the acceleration should be retarded by an amount equal to r/c , if the retardations are not equal, the signals are no longer in phase. Thus it should be possible to find a position at which the distances of D from S_1 and S_2 differ by some amount Δ , in such a manner that there is no net signal. That is, the distance Δ is to be the distance light goes in one-half an oscillation of the generator. We may go still further, and find a point where the difference is greater by a whole cycle; that is to say, the signal from the first antenna reaches point 3 with a delay in time that is greater than that of the second antenna by just the length of time it takes for the electric current to oscillate once, and therefore the two electric fields produced at 3 are in phase again. At point 3 the signal is strong again.

This completes our discussion of the experimental verification of some of the important features of Eq. (28.6). Of course we have not really checked the $1/r$ variation of the electric field strength, or the fact that there is also a magnetic field that goes along with the electric field. To do so would require rather sophisticated techniques and would hardly add to our understanding at this point. In any case, we have checked those features that are of the greatest importance for our later applications, and we shall come back to study some of the other properties of electromagnetic waves next year.

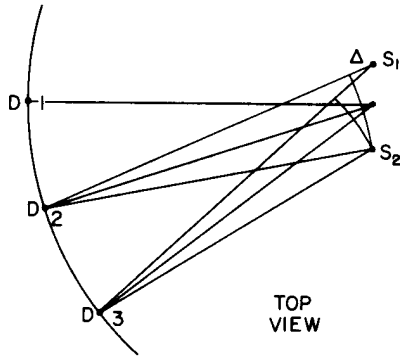


Fig. 28-3. Illustration of interference of sources.

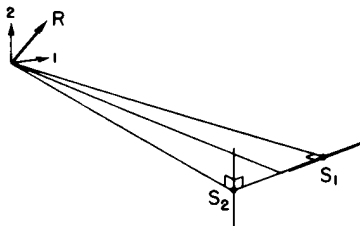


Fig. 28-4. Illustration of the vector character of the combination of sources.

Interference

29-1 Electromagnetic waves

In this chapter we shall discuss the subject of the preceding chapter more mathematically. We have qualitatively demonstrated that there are maxima and minima in the radiation field from two sources, and our problem now is to describe the field in mathematical detail, not just qualitatively.

We have already physically analyzed the meaning of formula (28.6) quite satisfactorily, but there are a few points to be made about it mathematically. In the first place, if a charge is accelerating up and down along a line, in a motion of very small amplitude, the field at some angle θ from the axis of the motion is in a direction at right angles to the line of sight and in the plane containing both the acceleration and the line of sight (Fig. 29-1). If the distance is called r , then at time t the electric field has the magnitude

$$E(t) = \frac{-qa(t - r/c) \sin \theta}{4\pi\epsilon_0 c^2 r}, \quad (29.1)$$

where $a(t - r/c)$ is the acceleration at the time $(t - r/c)$, called the *retarded* acceleration.

Now it would be interesting to draw a picture of the field under different conditions. The thing that is interesting, of course, is the factor $a(t - r/c)$, and to understand it we can take the simplest case, $\theta = 90^\circ$, and plot the field graphically. What we had been thinking of before is that we stand in one position and ask how the field there changes with time. But instead of that, we are now going to see what the field looks like at different positions in space at a given instant. So what we want is a "snapshot" picture which tells us what the field is in different places. Of course it depends upon the acceleration of the charge. Suppose that the charge at first had some particular motion: it was initially standing still, and it suddenly accelerated in some manner, as shown in Fig. 29-2, and then stopped. Then, a little bit later, we measure the field at a different place. Then we may assert that the field will appear as shown in Fig. 29-3. At each point the field is determined by the acceleration of the charge at an earlier time, the amount earlier being the delay r/c . The field at farther and farther points is determined by the acceleration at earlier and earlier times. So the curve in Fig. 29-3 is really, in a sense, a "reversed" plot of the acceleration as a function of time; the distance is related to time by a constant scale factor c , which we often take as unity. This is easily seen by considering the mathematical behavior of $a(t - r/c)$. Evidently, if we add a little time Δt , we get the same value for $a(t - r/c)$ as we would have if we had subtracted a little distance: $\Delta r = -c \Delta t$.

Stated another way: if we add a little time Δt , we can restore $a(t - r/c)$ to its former value by *adding* a little distance $\Delta r = c \Delta t$. That is, as time goes on *the field moves as a wave outward from the source*. That is the reason why we sometimes say light is propagated as waves. It is equivalent to saying that the field is delayed, or to saying that the electric field is moving outward as time goes on.

An interesting special case is that where the charge q is moving up and down in an oscillatory manner. The case which we studied experimentally in the last chapter was one in which the displacement x at any time t was equal to a certain constant x_0 , the magnitude of the oscillation, times $\cos \omega t$. Then the acceleration is

$$a = -\omega^2 x_0 \cos \omega t = a_0 \cos \omega t, \quad (29.2)$$

29-1 Electromagnetic waves

29-2 Energy of radiation

29-3 Sinusoidal waves

29-4 Two dipole radiators

29-5 The mathematics of interference

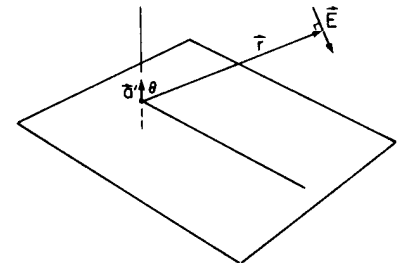


Fig. 29-1. The electric field \mathbf{E} due to a positive charge whose retarded acceleration is a' .

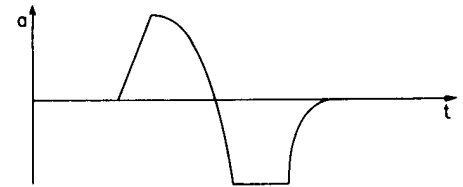


Fig. 29-2. The acceleration of a certain charge as a function of time.

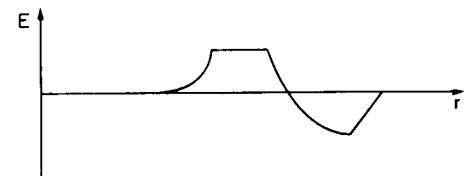


Fig. 29-3. The electric field as a function of position at a later time. (The $1/r$ variation is ignored.)

where a_0 is the maximum acceleration, $-\omega^2 x_0$. Putting this formula into (29.1), we find

$$E = -q \sin \theta \frac{a_0 \cos \omega(t - r/c)}{4\pi\epsilon_0 r c^2}. \quad (29.3)$$

Now, ignoring the angle θ and the constant factors, let us see what that looks like as a function of position or as a function of time.

29-2 Energy of radiation

First of all, at any particular moment or in any particular place, the strength of the field varies inversely as the distance r , as we mentioned previously. Now we must point out that the *energy* content of a wave, or the energy effects that such an electric field can have, are proportional to the *square* of the field, because if, for instance, we have some kind of a charge or an oscillator in the electric field, then if we let the field act on the oscillator, it makes it move. If this is a linear oscillator, the acceleration, velocity, and displacement produced by the electric field acting on the charge are all proportional to the field. So the kinetic energy which is developed in the charge is proportional to the *square* of the field. So we shall take it that the energy that a field can deliver to a system is proportional somehow to the square of the field.

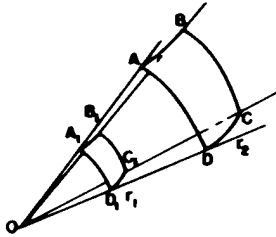


Fig. 29-4. The energy flowing within the cone OABCD is independent of the distance r at which it is measured.

This means that the energy that the source can deliver decreases as we get farther away; in fact, it varies *inversely as the square of the distance*. But that has a very simple interpretation: if we wanted to pick up all the energy we could from the wave in a certain cone at a distance r_1 (Fig. 29-4), and we do the same at another distance r_2 , we find that the amount of energy per unit area at any one place goes inversely as the square of r , but the area of the surface intercepted by the cone goes *directly* as the square of r . So the energy that we can take out of the wave within a given conical angle is the same, no matter how far away we are! In particular, the total energy that we could take out of the whole wave by putting absorbing oscillators all around is a certain fixed amount. So the fact that the amplitude of E varies as $1/r$ is the same as saying that there is an energy flux which is never lost, an energy which goes on and on, spreading over a greater and greater effective area. Thus we see that after a charge has oscillated, it has lost some energy which it can never recover; the energy keeps going farther and farther away without diminution. So if we are far enough away that our basic approximation is good enough, the charge cannot recover the energy which has been, as we say, radiated away. Of course the energy still exists somewhere, and is available to be picked up by other systems. We shall study this energy "loss" further in Chapter 32.

Let us now consider more carefully how the wave (29.3) varies as a function of time at a given place, and as a function of position at a given time. Again we ignore the $1/r$ variation and the constants.

29-3 Sinusoidal waves

First let us fix the position r , and watch the field as a function of time. It is oscillatory at the angular frequency ω . The angular frequency ω can be defined as the *rate of change of phase with time* (radians per second). We have already studied such a thing, so it should be quite familiar to us by now. The *period* is the time needed for one oscillation, one complete cycle, and we have worked that out too; it is $2\pi/\omega$, because ω times the period is one cycle of the cosine.

Now we introduce a new quantity which is used a great deal in physics. This has to do with the opposite situation, in which we fix t and look at the wave as a function of distance r . Of course we notice that, as a function of r , the wave (29.3) is also oscillatory. That is, aside from $1/r$, which we are ignoring, we see that E oscillates as we change the position. So, in analogy with ω , we can define a quantity called the *wave number*, symbolized as k . This is defined as the *rate of change of phase with distance* (radians per meter). That is, as we move in space at a fixed time, the phase changes.

There is another quantity that corresponds to the period, and we might call it the period in space, but it is usually called the wavelength, symbolized λ . The wavelength is the distance occupied by one complete cycle. It is easy to see, then, that the wavelength is $2\pi/k$, because k times the wavelength would be the number of radians that the whole thing changes, being the product of the rate of change of the radians per meter, times the number of meters, and we must make a 2π change for one cycle. So $k\lambda = 2\pi$ is exactly analogous to $\omega t_0 = 2\pi$.

Now in our particular wave there is a definite relationship between the frequency and the wavelength, but the above definitions of k and ω are actually quite general. That is, the wavelength and the frequency may not be related in the same way in other physical circumstances. However, in our circumstance the rate of change of phase with distance is easily determined, because if we call $\phi = \omega(t - r/c)$ the phase, and differentiate (partially) with respect to distance r , the rate of change, $\partial\phi/\partial r$, is

$$\left| \frac{\partial\phi}{\partial r} \right| = k = \frac{\omega}{c}. \quad (29.4)$$

There are many ways to represent the same thing, such as

$$\lambda = ct_0 \quad (29.5) \quad \lambda\nu = c \quad (29.7)$$

$$\omega = ck \quad (29.6) \quad \omega\lambda = 2\pi c \quad (29.8)$$

Why is the wavelength equal to c times the period? That's very easy, of course, because if we sit still and wait for one period to elapse, the waves, travelling at the speed c , will move a distance ct_0 , and will of course have moved over just one wavelength.

In a physical situation other than that of light, k is not necessarily related to ω in this simple way. If we call the distance along an axis x , then the formula for a cosine wave moving in a direction x with a wave number k and an angular frequency ω will be written in general as $\cos(\omega t - kx)$.

Now that we have introduced the idea of wavelength, we may say something more about the circumstances in which (29.1) is a legitimate formula. We recall that the field is made up of several pieces, one of which varies inversely as r , another part which varies inversely as r^2 , and others which vary even faster. It would be worth while to know in what circumstances the $1/r$ part of the field is the most important part, and the other parts are relatively small. Naturally, the answer is "if we go 'far enough' away," because terms which vary inversely as the square ultimately become negligible compared with the $1/r$ term. How far is "far enough"? The answer is, qualitatively, that the other terms are of order λ/r smaller than the $1/r$ term. Thus, so long as we are beyond a few wavelengths, (29.1) is an excellent approximation to the field. Sometimes the region beyond a few wavelengths is called the "wave zone."

29-4 Two dipole radiators

Next let us discuss the mathematics involved in combining the effects of two oscillators to find the net field at a given point. This is very easy in the few cases that we considered in the previous chapter. We shall first describe the effects qualitatively, and then more quantitatively. Let us take the simple case, where the oscillators are situated with their centers in the same horizontal plane as the detector, and the line of vibration is vertical.

Figure 29-5(a) represents the top view of two such oscillators, and in this particular example they are half a wavelength apart in a N-S direction, and are oscillating together in the same phase, which we call zero phase. Now we would like to know the intensity of the radiation in various directions. By the intensity we mean the amount of energy that the field carries past us per second, which is proportional to the square of the field, averaged in time. So the thing to look at, when we want to know how bright the light is, is the square of the electric field, not the electric field itself. (The electric field tells the strength of the force felt by a

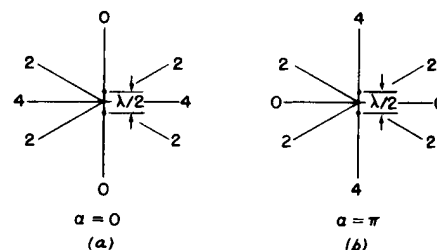


Fig. 29-5. The intensities in various directions from two dipole oscillators one-half wavelength apart. Left: in phase ($\alpha = 0$). Right: one-half period out of phase ($\alpha = \pi$).

stationary charge, but the amount of energy that is going past, in watts per square meter, is proportional to the square of the electric field. We shall derive the constant of proportionality in the next chapter.) If we look at the array from the W side, both oscillators contribute equally and in phase, so the electric field is twice as strong as it would be from a single oscillator. Therefore *the intensity is four times as strong as it would be if there were only one oscillator.* (The numbers in Fig. 29-5 represent how strong the intensity would be in this case, compared with what it would be if there were only a single oscillator of unit strength.) Now, in either the N or S direction along the line of the oscillators, since they are half a wavelength apart, the effect of one oscillator turns out to be out of phase by exactly half an oscillation from the other, and therefore the fields add to zero. At a certain particular intermediate angle (in fact, at 30°) the intensity is 2, and it falls off, 4, 2, 0, and so forth. We have to learn how to find these numbers at other angles. It is a question of adding two oscillations with different phases.

Let us quickly look at some other cases of interest. Suppose the oscillators are again one-half a wavelength apart, but the phase α of one is set half a period behind the other in its oscillation (Fig. 29-5b). In the W direction the intensity is now zero, because one oscillator is “pushing” when the other one is “pulling.” But in the N direction the signal from the near one comes at a certain time, and that of the other comes half a period later. But the latter was *originally* half a period behind in timing, and therefore it is now exactly *in time* with the first one, and so the intensity in this direction is 4 units. The intensity in the direction at 30° is still 2, as we can prove later.

Now we come to an interesting case which shows up a possibly useful feature. Let us remark that one of the reasons that phase relations of oscillators are interesting is for beaming radio transmitters. For instance, if we build an antenna system and want to send a radio signal, say, to Hawaii, we set the antennas up as in Fig. 29-5(a) and we broadcast with our two antennas in phase, because Hawaii is to the west of us. Then we decide that tomorrow we are going to broadcast toward Alberta, Canada. Since that is north, not west, all we have to do is to reverse the phase of one of our antennas, and we can broadcast to the north. So we can build antenna systems with various arrangements. Ours is one of the simplest possible ones; we can make them much more complicated, and by changing the phases in the various antennas we can send the beams in various directions and send most of the power in the direction in which we wish to transmit, without ever moving the antenna! In both of the preceding cases, however, while we are broadcasting toward Alberta we are wasting a lot of power on Easter Island, and it would be interesting to ask whether it is possible to send it in only *one* direction. At first sight we might think that with a pair of antennas of this nature the result is always going to be symmetrical. So let us consider a case that comes out unsymmetrical, to show the possible variety.

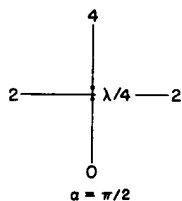


Fig. 29-6. A pair of dipole antennas giving maximum power in one direction.

If the antennas are separated by one-quarter wavelength, and if the N one is one-fourth period behind the S one in time, then what happens (Fig. 29-6)? In the W direction we get 2, as we will see later. In the S direction we get *zero*, because the signal from S comes at a certain time; that from N comes 90° later in *time*, but it is already 90° behind in its built-in phase, therefore it arrives, altogether, 180° out of phase, and there is no effect. On the other hand, in the N direction, the N signal arrives earlier than the S signal by 90° in time, because it is a quarter wavelength closer. But its phase is set so that it is oscillating 90° *behind* in time, which just compensates the delay difference, and therefore the two signals appear *together* in phase, making the field strength twice as large, and the energy four times as great.

Thus, by using some cleverness in spacing and phasing our antennas, we can send the power all in one direction. But still it is distributed over a great range of angles. Can we arrange it so that it is focused still more sharply in a particular direction? Let us consider the case of Hawaii again, where we are sending the beam east and west but it is spread over quite an angle, because even at 30° we are still getting half the intensity—we are wasting the power. Can we do better than that? Let us take a situation in which the separation is ten wavelengths (Fig.

29–7), which is more nearly comparable to the situation in which we experimented in the previous chapter, with separations of several wavelengths rather than a small fraction of a wavelength. Here the picture is quite different.

If the oscillators are ten wavelengths apart (we take the in-phase case to make it easy), we see that in the E–W direction, they are in phase, and we get a strong intensity, four times what we would get if one of them were there alone. On the other hand, at a very small angle away, the arrival times differ by 180° and the intensity is zero. To be precise, if we draw a line from each oscillator to a distant point and the difference Δ in the two distances is $\lambda/2$, half an oscillation, then they will be out of phase. So this first null occurs when that happens. (The figure is not drawn to scale; it is only a rough sketch.) This means that we do indeed have a very sharp beam in the direction we want, because if we just move over a little bit we lose all our intensity. Unfortunately for practical purposes, if we were thinking of making a radio broadcasting array and we doubled the distance Δ , then we would be a whole cycle out of phase, which is the same as being exactly *in* phase again! Thus we get many successive maxima and minima, just as we found with the $2\frac{1}{2}\lambda$ spacing in Chapter 28.

Now how can we arrange to get rid of all these extra maxima, or “lobes,” as they are called? We could get rid of the unwanted lobes in a rather interesting way. Suppose that we were to place another set of antennas between the two that we already have. That is, the outside ones are still 10λ apart, but between them, say every 2λ , we have put another antenna, and we drive them all in phase. There are now six antennas, and if we looked at the intensity in the E–W direction, it would, of course, be much higher with six antennas than with one. The field would be six times and the intensity thirty-six times as great (the square of the field). We get 36 units of intensity in that direction. Now if we look at neighboring points, we find a zero as before, roughly, but if we go farther, to where we used to get a big “bump,” we get a much smaller “bump” now. Let us try to see why.

The reason is that although we might expect to get a big bump when the distance Δ is exactly equal to the wavelength, it is true that dipoles 1 and 6 are then in phase and are cooperating in trying to get some strength in that direction. But numbers 3 and 4 are roughly $\frac{1}{2}$ a wavelength out of phase with 1 and 6, and although 1 and 6 push together, 3 and 4 push together too, but in opposite phase. Therefore there is very little intensity in this direction—but there is something; it does not balance exactly. This kind of thing keeps on happening; we get very little bumps, and we have the strong beam in the direction where we want it. But in this particular example, something else will happen: namely, since the distance between successive dipoles is 2λ , it is possible to find an angle where the distance δ between successive dipoles is exactly one wavelength, so that the effects from all of them are in phase again. Each one is delayed relative to the next one by 360° , so they all come back in phase, and we have another strong beam in that direction! It is easy to avoid this in practice because it is possible to put the dipoles closer than one wavelength apart. If we put in more antennas, closer than one wavelength apart, then this cannot happen. But the fact that this *can* happen at certain angles, if the spacing is bigger than one wavelength, is a very interesting and useful phenomenon in other applications—not radio broadcasting, but in *diffraction gratings*.

29–5 The mathematics of interference

Now we have finished our analysis of the phenomena of dipole radiators qualitatively, and we must learn how to analyze them quantitatively. To find the effect of two sources at some particular angle in the most general case, where the two oscillators have some intrinsic relative phase α from one another and the strengths A_1 and A_2 are not equal, we find that we have to add two cosines having the same frequency, but with different phases. It is very easy to find this phase difference; it is made up of a delay due to the difference in distance, and the intrinsic, built-in phase of the oscillation. Mathematically, we have to find the sum R of two waves: $R = A_1 \cos(\omega t + \phi_1) + A_2 \cos(\omega t + \phi_2)$. How do we do it?

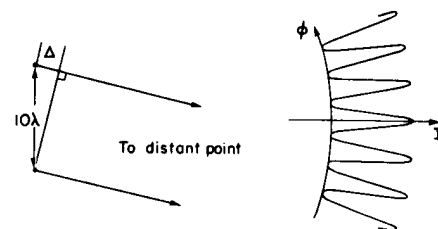


Fig. 29–7. The intensity pattern for two dipoles separated by 10λ .

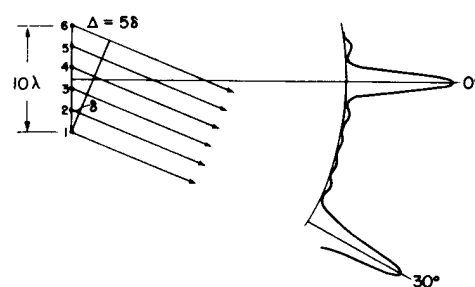


Fig. 29–8. A six-dipole antenna array and part of its intensity pattern.

It is really very easy, and we presume that we already know how to do it. However, we shall outline the procedure in some detail. First, we can, if we are clever with mathematics and know enough about cosines and sines, simply work it out. The easiest such case is the one where A_1 and A_2 are equal, let us say they are both equal to A . In those circumstances, for example (we could call this the trigonometric method of solving the problem), we have

$$R = A[\cos(\omega t + \phi_1) + \cos(\omega t + \phi_2)]. \quad (29.9)$$

Once, in our trigonometry class, we may have learned the rule that

$$\cos A + \cos B = 2 \cos \frac{1}{2}(A + B) \cos \frac{1}{2}(A - B). \quad (29.10)$$

If we know that, then we can immediately write R as

$$R = 2A \cos \frac{1}{2}(\phi_1 - \phi_2) \cos(\omega t + \frac{1}{2}\phi_1 + \frac{1}{2}\phi_2). \quad (29.11)$$

So we find that we have an oscillatory wave with a new phase and a new amplitude. In general, the result *will* be an oscillatory wave with a new amplitude A_R , which we may call the resultant amplitude, oscillating at the same frequency but with a phase difference ϕ_R , called the resultant phase. In view of this, our particular case has the following result: that the resultant amplitude is

$$A_R = 2A \cos \frac{1}{2}(\phi_1 - \phi_2), \quad (29.12)$$

and the resultant phase is the average of the two phases, and we have completely solved our problem.

Now suppose that we cannot remember that the sum of two cosines is twice the cosine of half the sum times the cosine of half the difference. Then we may use another method of analysis which is more geometrical. Any cosine function of ωt can be considered as the horizontal projection of a *rotating vector*. Suppose there were a vector \mathbf{A}_1 of length A_1 rotating with time, so that its angle with the horizontal axis is $\omega t + \phi_1$. (We shall leave out the ωt in a minute, and see that it makes no difference.) Suppose that we take a snapshot at the time $t = 0$, although, in fact, the picture is rotating with angular velocity ω (Fig. 29-9). The projection of \mathbf{A}_1 along the horizontal axis is precisely $A_1 \cos(\omega t + \phi_1)$. Now at $t = 0$ the second wave could be represented by another vector, \mathbf{A}_2 , of length A_2 and at an angle ϕ_2 , and also rotating. They are both rotating with the same angular velocity ω , and therefore the *relative* positions of the two are fixed. The system goes around like a rigid body. The horizontal projection of \mathbf{A}_2 is $A_2 \cos(\omega t + \phi_2)$. But we know from the theory of vectors that if we add the two vectors in the ordinary way, by the parallelogram rule, and draw the resultant vector \mathbf{A}_R , the x-component of the resultant is the sum of the x-components of the other two vectors. That solves our problem. It is easy to check that this gives the correct result for the special case we treated above, where $A_1 = A_2 = A$. In this case, we see from Fig. 29-9 that \mathbf{A}_R lies midway between \mathbf{A}_1 and \mathbf{A}_2 and makes an angle $\frac{1}{2}(\phi_2 - \phi_1)$ with each. Therefore we see that $A_R = 2A \cos \frac{1}{2}(\phi_2 - \phi_1)$, as before. Also, as we see from the triangle, the phase of \mathbf{A}_R , as it goes around, is the average angle of \mathbf{A}_1 and \mathbf{A}_2 when the two amplitudes are equal. Clearly, we can also solve for the case where the amplitudes are not equal, just as easily. We can call that the *geometrical* way of solving the problem.

There is still another way of solving the problem, and that is the *analytical* way. That is, instead of having actually to draw a picture like Fig. 29-9, we can write something down which says the same thing as the picture: instead of drawing the vectors, we write a *complex number* to represent each of the vectors. The real parts of the complex numbers are the actual physical quantities. So in our particular case the waves could be written in this way: $A_1 e^{i(\omega t + \phi_1)}$ [the real part of this is $A_1 \cos(\omega t + \phi_1)$] and $A_2 e^{i(\omega t + \phi_2)}$. Now we can add the two:

$$R = A_1 e^{i(\omega t + \phi_1)} + A_2 e^{i(\omega t + \phi_2)} = (A_1 e^{i\phi_1} + A_2 e^{i\phi_2}) e^{i\omega t} \quad (29.13)$$

or

$$\hat{R} = A_1 e^{i\phi_1} + A_2 e^{i\phi_2} = A_R e^{i\phi_R}. \quad (29.14)$$

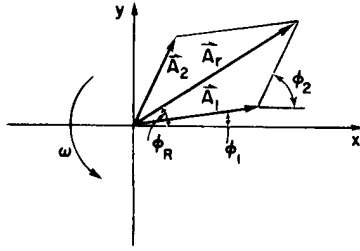


Fig. 29-9. A geometrical method for combining two cosine waves. The entire diagram is thought of as rotating counter-clockwise with angular frequency ω .

This solves the problem that we wanted to solve, because it represents the result as a complex number of magnitude A_R and phase ϕ_R .

To see how this method works, let us find the amplitude A_R which is the “length” of \hat{R} . To get the “length” of a complex quantity, we always multiply the quantity by its complex conjugate, which gives the length squared. The complex conjugate is the same expression, but with the sign of the i ’s reversed. Thus we have

$$A_R^2 = (A_1 e^{i\phi_1} + A_2 e^{i\phi_2})(A_1 e^{-i\phi_1} + A_2 e^{-i\phi_2}). \quad (29.15)$$

In multiplying this out, we get $A_1^2 + A_2^2$ (here the e ’s cancel), and for the cross terms we have

$$A_1 A_2 (e^{i(\phi_1 - \phi_2)} + e^{i(\phi_2 - \phi_1)}).$$

Now

$$e^{i\theta} + e^{-i\theta} = \cos \theta + i \sin \theta + \cos \theta - i \sin \theta.$$

That is to say, $e^{i\theta} + e^{-i\theta} = 2 \cos \theta$. Our final result is therefore

$$A_R^2 = A_1^2 + A_2^2 + 2A_1 A_2 \cos(\phi_2 - \phi_1). \quad (29.16)$$

As we see, this agrees with the length of A_R in Fig. 29-9, using the rules of trigonometry.

Thus the sum of the two effects has the intensity A_1^2 we would get with one of them alone, plus the intensity A_2^2 we would get with the other one alone, plus a correction. This correction we call the *interference effect*. It is really only the difference between what we get simply by adding the intensities, and what actually happens. We call it interference whether it is positive or negative. (Interference in ordinary language usually suggests opposition or hindrance, but in physics we often do not use language the way it was originally designed!) If the interference term is positive, we call that case *constructive* interference, horrible though it may sound to anybody other than a physicist! The opposite case is called *destructive* interference.

Now let us see how to apply our general formula (29.16) for the case of two oscillators to the special situations which we have discussed qualitatively. To apply this general formula, it is only necessary to find what phase difference, $\phi_1 - \phi_2$, exists between the signals arriving at a given point. (It depends only on the phase difference, of course, and not on the phase itself.) So let us consider the case where the two oscillators, of equal amplitude, are separated by some distance d and have an intrinsic relative phase α . (When one is at phase zero, the phase of the other is α .) Then we ask what the intensity will be in some azimuth direction θ from the E-W line. [Note that this is *not* the same θ as appears in (29.1). We are torn between using an unconventional symbol like \mathcal{U} , or the conventional symbol θ (Fig. 29-10).] The phase relationship is found by noting that the difference in distance from P to the two oscillators is $d \sin \theta$, so that the phase difference contribution from this is the number of wavelengths in $d \sin \theta$, multiplied by 2π . (Those who are more sophisticated might want to multiply the wave number k , which is the rate of change of phase with distance, by $d \sin \theta$; it is exactly the same.) The phase difference due to the distance difference is thus $2\pi d \sin \theta / \lambda$, but, due to the timing of the oscillators, there is an additional phase α . So the phase difference at arrival would be

$$\phi_2 - \phi_1 = \alpha + 2\pi d \sin \theta / \lambda. \quad (29.17)$$

This takes care of all the cases. Thus all we have to do is substitute this expression into (29.16) for the case $A_1 = A_2$, and we can calculate all the various results for two antennas of equal intensity.

Now let us see what happens in our various cases. The reason we know, for example, that the intensity is 2 at 30° in Fig. 29-5 is the following: the two oscillators are $\frac{1}{2}\lambda$ apart, so at 30° , $d \sin \theta = \lambda/4$. Thus $\phi_2 - \phi_1 = 2\pi\lambda/4\lambda = \pi/2$, and so the interference term is zero. (We are adding two vectors at 90° .) The result is the hypotenuse of a 45° right-angle triangle, which is $\sqrt{2}$ times the unit amplitude; squaring it, we get twice the intensity of one oscillator alone. All the other cases can be worked out in this same way.

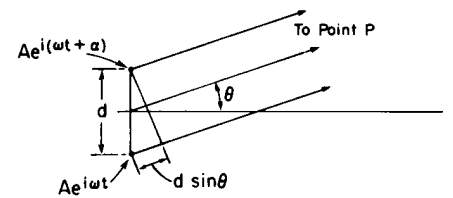


Fig. 29-10. Two oscillators of equal amplitude, with a phase difference α between them.

Diffraction

30-1 The resultant amplitude due to n equal oscillators

This chapter is a direct continuation of the previous one, although the name has been changed from *Interference* to *Diffraction*. No one has ever been able to define the difference between interference and diffraction satisfactorily. It is just a question of usage, and there is no specific, important physical difference between them. The best we can do, roughly speaking, is to say that when there are only a few sources, say two, interfering, then the result is usually called interference, but if there is a large number of them, it seems that the word diffraction is more often used. So, we shall not worry about whether it is interference or diffraction, but continue directly from where we left off in the middle of the subject in the last chapter.

Thus we shall now discuss the situation where there are n equally spaced oscillators, all of equal amplitude but different from one another in phase, either because they are driven differently in phase, or because we are looking at them at an angle such that there is a difference in time delay. For one reason or another, we have to add something like this:

$$R = A[\cos \omega t + \cos(\omega t + \phi) + \cos(\omega t + 2\phi) + \cdots + \cos(\omega t + (n-1)\phi)], \quad (30.1)$$

where ϕ is the phase difference between one oscillator and the next one, as seen in a particular direction. Specifically, $\phi = \alpha + 2\pi d \sin \theta / \lambda$. Now we must add all the terms together. We shall do this geometrically. The first one is of length A , and it has zero phase. The next is also of length A and it has a phase equal to ϕ . The next one is again of length A and it has a phase equal to 2ϕ , and so on. So we are evidently going around an equiangular polygon with n sides (Fig. 30-1).

Now the vertices, of course, all lie on a circle, and we can find the net amplitude most easily if we find the radius of that circle. Suppose that Q is the center of the circle. Then we know that the angle OQS is just a phase angle ϕ . (This is because the radius QS bears the same geometrical relation to A_2 as QO bears to A_1 , so they form an angle ϕ between them.) Therefore the radius r must be such that $A = 2r \sin \phi/2$, which fixes r . But the large angle OQT is equal to $n\phi$, and we thus find that $A_R = 2r \sin n\phi/2$. Combining these two results to eliminate r , we get

$$A_R = A \frac{\sin n\phi/2}{\sin \phi/2}. \quad (30.2)$$

The resultant intensity is thus

$$I = I_0 \frac{\sin^2 n\phi/2}{\sin^2 \phi/2}. \quad (30.3)$$

Now let us analyze this expression and study some of its consequences. In the first place, we can check it for $n = 1$. It checks: $I = I_0$. Next, we check it for $n = 2$: writing $\sin \phi = 2 \sin \phi/2 \cos \phi/2$, we find that $A_R = 2A \cos \phi/2$, which agrees with (29.12).

Now the idea that led us to consider the addition of several sources was that we might get a much stronger intensity in one direction than in another; that the nearby maxima which would have been present if there were only two sources will have gone down in strength. In order to see this effect, we plot the curve that comes from (30.3), taking n to be enormously large and plotting the region near

30-1 The resultant amplitude due to n equal oscillators

30-2 The diffraction grating

30-3 Resolving power of a grating

30-4 The parabolic antenna

30-5 Colored films; crystals

30-6 Diffraction by opaque screens

30-7 The field of a plane of oscillating charges

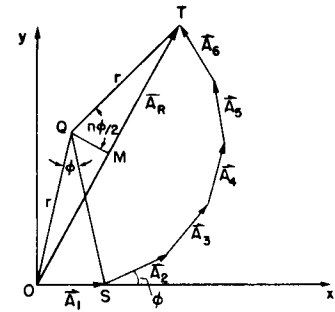


Fig. 30-1. The resultant amplitude of $n = 6$ equally spaced sources with net successive phase differences ϕ .

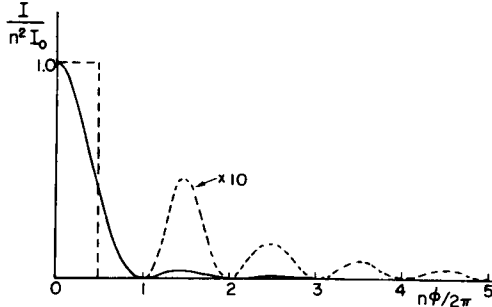


Fig. 30-2. The intensity as a function of phase angle for a large number of oscillators of equal strength.

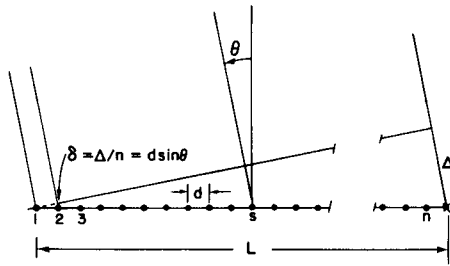


Fig. 30-3. A linear array of n equal oscillators, driven with phases $\alpha_s = s\alpha$.

$\phi = 0$. In the first place, if ϕ is exactly 0, we have 0/0, but if ϕ is infinitesimal, the ratio of the two sines squared is simply n^2 , since the sine and the angle are approximately equal. Thus the intensity of the maximum of the curve is equal to n^2 times the intensity of one oscillator. That is easy to see, because if they are all in phase, then the little vectors have no relative angle and all n of them add up so the amplitude is n times, and the intensity n^2 times, stronger.

As the phase ϕ increases, the ratio of the two sines begins to fall off, and the first time it reaches zero is when $n\phi/2 = \pi$, because $\sin \pi = 0$. In other words, $\phi = 2\pi/n$ corresponds to the first minimum in the curve (Fig. 30-2). In terms of what is happening with the arrows in Fig. 30-1, the first minimum occurs when all the arrows come back to the starting point; that means that the total accumulated angle in all the arrows, the total phase difference between the first and last oscillator, must be 2π to complete the circle.

Now we go to the next maximum, and we want to see that it is really much smaller than the first one, as we had hoped. We shall not go precisely to the maximum position, because both the numerator and the denominator of (30.3) are variant, but $\sin \phi/2$ varies quite slowly compared with $\sin n\phi/2$ when n is large, so when $\sin n\phi/2 = 1$ we are very close to the maximum. The next maximum of $\sin^2 n\phi/2$ comes at $n\phi/2 = 3\pi/2$, or $\phi = 3\pi/n$. This corresponds to the arrows having traversed the circle one and a half times. On putting $\phi = 3\pi/n$ into the formula to find the size of the maximum, we find that $\sin^2 3\pi/2 = 1$ in the numerator (because that is why we picked this angle), and in the denominator we have $\sin^2 3\pi/2n$. Now if n is sufficiently large, then this angle is very small and the sine is equal to the angle; so for all practical purposes, we can put $\sin 3\pi/2n = 3\pi/2n$. Thus we find that the intensity at this maximum is $I = I_0(4n^2/9\pi^2)$. But $n^2 I_0$ was the maximum intensity, and so we have $4/9\pi^2$ times the maximum intensity, which is about 0.047, less than 5 percent, of the maximum intensity! Of course there are decreasing intensities farther out. So we have a very sharp central maximum with very weak subsidiary maxima on the sides.

It is possible to prove that the area of the whole curve, including all the little bumps, is equal to $2\pi n I_0$, or twice the area of the dotted rectangle in Fig. 30-2.

Now let us consider further how we may apply Eq. (30.3) in different circumstances, and try to understand what is happening. Let us consider our sources to be all on a line, as drawn in Fig. 30-3. There are n of them, all spaced by a distance d , and we shall suppose that the intrinsic relative phase, one to the next, is α . Then if we are observing in a given direction θ from the normal, there is an additional phase $2\pi d \sin \theta / \lambda$ because of the time delay between each successive two, which we talked about before. Thus

$$\begin{aligned}\phi &= \alpha + 2\pi d \sin \theta / \lambda \\ &= \alpha + kd \sin \theta.\end{aligned}\quad (30.4)$$

First, we shall take the case $\alpha = 0$. That is, all oscillators are in phase, and we want to know what the intensity is as a function of the angle θ . In order to find out, we merely have to put $\phi = kd \sin \theta$ into formula (30.3) and see what happens. In the first place, there is a maximum when $\phi = 0$. That means that when all the oscillators are in phase there is a strong intensity in the direction $\theta = 0$. On the other hand, an interesting question is, where is the first minimum? That occurs when $\phi = 2\pi/n$. In other words, when $2\pi d \sin \theta / \lambda = 2\pi/n$, we get the first minimum of the curve. If we get rid of the 2π 's so we can look at it a little better, it says that

$$nd \sin \theta = \lambda. \quad (30.5)$$

Now let us understand physically why we get a minimum at that position. nd is the total length L of the array. Referring to Fig. 30-3, we see that $nd \sin \theta = L \sin \theta = \Delta$. What (30.5) says is that when Δ is equal to one wavelength, we get a minimum. Now why do we get a minimum when $\Delta = \lambda$? Because the contributions of the various oscillators are then uniformly distributed in phase from 0° to

360°. The arrows (Fig. 30-1) are going around a whole circle—we are adding equal vectors in all directions, and such a sum is zero. So when we have an angle such that $\Delta = \lambda$, we get a minimum. That is the first minimum.

There is another important feature about formula (30.3), which is that if the angle ϕ is increased by any multiple of 2π , it makes no difference to the formula. So we will get other strong maxima at $\phi = 2\pi, 4\pi, 6\pi$, and so forth. Near each of these great maxima the pattern of Fig. 30-2 is repeated. We may ask ourselves, what is the geometrical circumstance that leads to these other great maxima? The condition is that $\phi = 2\pi m$, where m is any integer. That is, $2\pi d \sin \theta / \lambda = 2\pi m$. Dividing by 2π , we see that

$$d \sin \theta = m\lambda. \quad (30.6)$$

This looks like the other formula, (30.5). No, that formula was $nd \sin \theta = \lambda$. The difference is that here we have to look at the *individual sources*, and when we say $d \sin \theta = m\lambda$, that means that we have an angle θ such that $\delta = m\lambda$. In other words, each source is now contributing a certain amount, and successive ones are out of phase by a whole multiple of 360°, and therefore are contributing *in phase*, because out of phase by 360° is the same as being in phase. So they all contribute in phase and produce just as good a maximum as the one for $m = 0$ that we discussed before. The subsidiary bumps, the whole shape of the pattern, is just like the one near $\phi = 0$, with exactly the same minima on each side, etc. Thus such an array will send beams in various directions—each beam having a strong central maximum and a certain number of weak “side lobes.” The various strong beams are referred to as the zero-order beam, the first-order beam, etc., according to the value of m . m is called the *order* of the beam.

We call attention to the fact that if d is less than λ , Eq. (30.6) can have no solution except $m = 0$, so that if the spacing is too small there is only one possible beam, the zero-order one centered at $\theta = 0$. (Of course, there is also a beam in the opposite direction.) In order to get subsidiary great maxima, we must have the spacing d of the array greater than one wavelength.

30-2 The diffraction grating

In technical work with antennas and wires it is possible to arrange that all the phases of the little oscillators, or antennas, are equal. The question is whether and how we can do a similar thing with light. We cannot at the present time literally make little optical-frequency radio stations and hook them up with infinitesimal wires and drive them all with a given phase. But there is a very easy way to do what amounts to the same thing.

Suppose that we had a lot of parallel wires, equally spaced at a spacing d , and a radiofrequency source very far away, practically at infinity, which is generating an electric field which arrives at each one of the wires at the same phase (it is so far away that the time delay is the same for all of the wires). (One can work out cases with curved arrays, but let us take a plane one.) Then the external electric field will drive the electrons up and down in each wire. That is, the field which is coming from the original source will shake the electrons up and down, and in moving, these represent *new generators*. This phenomenon is called scattering: a light wave from some source can induce a motion of the electrons in a piece of material, and these motions generate their own waves. Therefore all that is necessary is to set up a lot of wires, equally spaced, drive them with a radiofrequency source far away, and we have the situation that we want, without a whole lot of special wiring. If the incidence is normal, the phases will be equal, and we will get exactly the circumstance we have been discussing. Therefore, if the wire spacing is greater than the wavelength, we will get a strong intensity of scattering in the normal direction, and in certain other directions given by (30.6).

This can also be done with light! Instead of wires, we use a flat piece of glass and make notches in it such that each of the notches scatters a little differently than the rest of the glass. If we then shine light on the glass, each one of the notches

will represent a source, and if we space the lines very finely, but not closer than a wavelength (which is technically almost impossible anyway), then we would expect a miraculous phenomenon: the light not only will pass straight through, but there will also be a strong beam at a finite angle, depending on the spacing of the notches! Such objects have actually been made and are in common use—they are called *diffraction gratings*.

In one of its forms, a diffraction grating consists of nothing but a plane glass sheet, transparent and colorless, with scratches on it. There are often several hundred scratches to the millimeter, *very* carefully arranged so as to be equally spaced. The effect of such a grating can be seen by arranging a projector so as to throw a narrow, vertical line of light (the image of a slit) onto a screen. When we put the grating into the beam, with its scratches vertical, we see that the line is still there but, in addition, on each side we have *another* strong patch of light which is *colored*. This, of course, is the slit image spread out over a wide angular range, because the angle θ in (30.6) depends upon λ , and lights of different colors, as we know, correspond to different frequencies, and therefore different wavelengths. The longest visible wavelength is red, and since $d \sin \theta = \lambda$, that requires a larger θ . And we do, in fact, find that red is at a greater angle out from the central image! There should also be a beam on the other side, and indeed we see one on the screen. Then, there might be another solution of (30.6) when $m = 2$. We do see that there is something vaguely there—very weak—and there are even other beams beyond.

We have just argued that all these beams ought to be of the same strength, but we see that they actually are not and, in fact, not even the first ones on the right and left are equal! The reason is that the grating has been carefully built to do just this. How? If the grating consists of very fine notches, infinitesimally wide, spaced evenly, then all the intensities would indeed be equal. But, as a matter of fact, although we have taken the simplest case, we could also have considered an array of *pairs* of antennas, in which each member of the pair has a certain strength and some relative phase. In this case, it is possible to get intensities which are different in the different orders. A grating is often made with little “sawtooth” cuts instead of little symmetrical notches. By carefully arranging the “sawteeth,” more light may be sent into one particular order of spectrum than into the others. In a practical grating, we would like to have as much light as possible in one of the orders. This may seem a complicated point to bring in, but it is a very clever thing to do, because it makes the grating more useful.

So far, we have taken the case where all the phases of the sources are equal. But we also have a formula for ϕ when the phases differ from one to the next by an angle α . That requires wiring up our antennas with a slight phase shift between each one. Can we do that with light? Yes, we can do it very easily, for suppose that there were a source of light at infinity, *at an angle* such that the light is coming in at an angle θ_{in} , and let us say that we wish to discuss the scattered beam, which is leaving at an angle θ_{out} . The θ_{out} is the same θ as we have had before, but the θ_{in} is merely a means for arranging that the phase of each source is different: the light coming from the distant driving source first hits one scratch, then the next, then the next, and so on, with a phase shift from one to the other, which, as we see, is $\alpha = -d \sin \theta_{in}/\lambda$. Therefore we have the formula for a grating in which light both comes in and goes out at an angle:

$$\phi = 2\pi d \sin \theta_{out}/\lambda - 2\pi d \sin \theta_{in}/\lambda. \quad (30.7)$$

Let us try to find out where we get strong intensity in these circumstances. The condition for strong intensities is, of course, that ϕ should be a multiple of 2π . There are several interesting points to be noted.

One case of rather great interest is that which corresponds to $m = 0$, where d is less than λ ; in fact, this is the only solution. In this case we see that $\sin \theta_{out} = \sin \theta_{in}$, which means that the light comes out in the *same direction* as the light which was exciting the grating. We might think that the light “goes right through.” No, it is *different light* that we are talking about. The light that goes right through is from the original source; what we are talking about is the new light *which is*

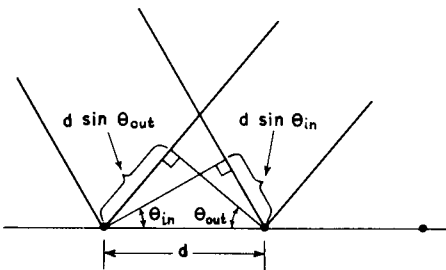


Fig. 30-4. The path difference for rays scattered from adjacent rulings of a grating is $d \sin \theta_{out} - d \sin \theta_{in}$.

generated by scattering. It turns out that the scattered light is going in the same direction as the original light, in fact it can interfere with it—a feature which we will study later.

There is another solution for this same case. For a given θ_{in} , θ_{out} may be the supplement of θ_{in} . So not only do we get a beam in the same direction as the incoming beam but also one in another direction, which, if we consider it carefully, is such that the *angle of incidence is equal to the angle of scattering*. This we call the *reflected beam*.

So we begin to understand the basic machinery of reflection: the light that comes in generates motions of the atoms in the reflector, and the reflector then regenerates a *new wave*, and one of the solutions for the direction of scattering, the *only* solution if the spacing of the scatterers is small compared with one wavelength, is that the angle at which the light comes out is equal to the angle at which it comes in!

Next, we discuss the special case when $d \rightarrow 0$. That is, we have just a solid piece of material, so to speak, but of finite length. In addition, we want the phase shift from one scatterer to the next to go to zero. In other words, we put more and more antennas between the other ones, so that each of the phase differences is getting smaller, but the number of antennas is increasing in such a way that the total phase difference, between one end of the line and the other, is constant. Let us see what happens to (30.3) if we keep the difference in phase $n\phi$ from one end to the other constant (say $n\phi = \Phi$), letting the number go to infinity and the phase shift ϕ of each one go to zero. But now ϕ is so small that $\sin \phi = \phi$, and if we also recognize $n^2 I_0$ as I_m , the maximum intensity at the center of the beam, we find

$$I = 4I_m \sin^2 \frac{1}{2}\Phi / \Phi^2. \quad (30.8)$$

This limiting case is what is shown in Fig. 30-2.

In such circumstances we find the same general kind of a picture as for finite spacing with $d > \lambda$; all the side lobes are practically the same as before, but there are no higher-order maxima. If the scatterers are all in phase, we get a maximum in the direction $\theta_{\text{out}} = 0$, and a minimum when the distance Δ is equal to λ , just as for finite d and n . So we can even analyze a *continuous* distribution of scatterers or oscillators, by using integrals instead of summing.

As an example, suppose there were a long line of oscillators, with the charge oscillating along the direction of the line (Fig. 30-5). From such an array the greatest intensity is perpendicular to the line. There is a little bit of intensity up and down from the equatorial plane, but it is very slight. With this result, we can handle a more complicated situation. Suppose we have a set of such lines, each producing a beam only in a plane perpendicular to the line. To find the intensity in various directions from a series of long wires, instead of infinitesimal wires, is the same problem as it was for infinitesimal wires, so long as we are in the central plane perpendicular to the wires; we just add the contribution from each of the long wires. That is why, although we actually analyzed only tiny antennas, we might as well have used a grating with long, narrow slots. Each of the long slots produces an effect only in its own direction, not up and down, but they are all set next to each other horizontally, so they produce interference that way.

Thus we can build up more complicated situations by having various distributions of scatterers in lines, planes, or in space. The first thing we did was to consider scatterers in a line, and we have just extended the analysis to strips; we can work it out by just doing the necessary summations, adding the contributions from the individual scatterers. The principle is always the same.

30-3 Resolving power of a grating

We are now in a position to understand a number of interesting phenomena. For example, consider the use of a grating for separating wavelengths. We noticed that the whole spectrum was spread out on the screen, so a grating can be used as an instrument for separating light into its different wavelengths. One of the

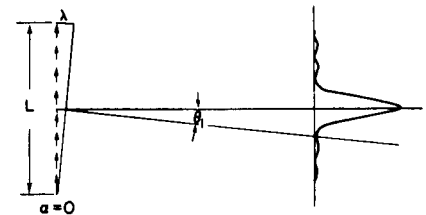


Fig. 30-5. The intensity pattern of a continuous line of oscillators has a single strong maximum and many weak "side lobes."

interesting questions is: supposing that there were two sources of slightly different frequency, or slightly different wavelength, how close together in wavelength could they be such that the grating would be unable to tell that there were really two different wavelengths there? The red and the blue were clearly separated. But when one wave is red and the other is slightly redder, very close, how close can they be? This is called the *resolving power* of the grating, and one way of analyzing the problem is as follows. Suppose that for light of a certain color we happen to have the maximum of the diffracted beam occurring at a certain angle. If we vary the wavelength the phase $2\pi d \sin \theta / \lambda$ is different, so of course the maximum occurs at a different angle. That is why the red and blue are spread out. How different in angle must it be in order for us to be able to see it? If the two maxima are exactly on top of each other, of course we cannot see them. If the maximum of one is far enough away from the other, then we can see that there is a double bump in the distribution of light. In order to be able to just make out the double bump, the following simple criterion, called *Rayleigh's criterion*, is usually used (Fig. 30-6). It is that the first minimum from one bump should sit at the maximum of the other. Now it is very easy to calculate, when one minimum sits on the other maximum, how much the difference in wavelength is. The best way to do it is geometrically.

In order to have a maximum for wavelength λ' , the distance Δ (Fig. 30-3) must be $n\lambda'$, and if we are looking at the m th-order beam, it is $mn\lambda'$. In other words, $2\pi d \sin \theta / \lambda' = 2\pi m$, so $nd \sin \theta$, which is Δ , is λ' times n , or $mn\lambda'$. For the other beam, of wavelength λ , we want to have a *minimum* at this angle. That is, we want Δ to be exactly one wavelength λ more than $mn\lambda$. That is, $\Delta = mn\lambda + \lambda = mn\lambda'$. Thus if $\lambda' = \lambda + \Delta\lambda$, we find

$$\Delta\lambda / \lambda = 1 / mn. \quad (30.9)$$

The ratio $\lambda / \Delta\lambda$ is called the *resolving power* of a grating; we see that it is equal to the total number of lines in the grating, times the order. It is not hard to prove that this formula is equivalent to the formula that the error in frequency is equal to the reciprocal time difference between extreme paths that are allowed to interfere.*

$$\Delta\nu = 1/T.$$

In fact, that is the best way to remember it, because the general formula works not only for gratings, but for any other instrument whatsoever, while the special formula (30.9) depends on the fact that we are using a grating.

30-4 The parabolic antenna

Now let us consider another problem in resolving power. This has to do with the antenna of a radio telescope, used for determining the position of radio sources in the sky, i.e., how large they are in angle. Of course if we use any old antenna and find signals, we would not know from what direction they came. We are very interested to know whether the source is in one place or another. One way we can find out is to lay out a whole series of equally spaced dipole wires on the Australian landscape. Then we take all the wires from these antennas and feed them into the same receiver, in such a way that all the delays in the feed lines are equal. Thus the receiver receives signals from all of the dipoles in phase. That is, it adds all the waves from every one of the dipoles in the same phase. Now what happens? If the source is directly above the array, at infinity or nearly so, then its radiowaves will excite all the antennas in the same phase, so they all feed the receiver together.

Now suppose that the radio source is at a slight angle θ from the vertical. Then the various antennas are receiving signals a little out of phase. The receiver adds all these out-of-phase signals together, and so we get nothing, if the angle

* In our case $T = \Delta / c = mn\lambda / c$, where c is the speed of light. The frequency $\nu = c / \lambda$, so $\Delta\nu = c \Delta\lambda / \lambda^2$.

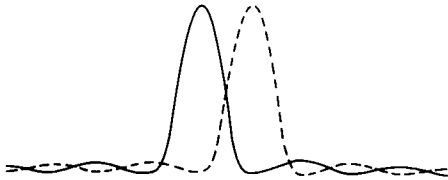


Fig. 30-6. Illustration of the Rayleigh criterion. The maximum of one pattern falls on the first minimum of the other.

θ is too big. How big may the angle be? *Answer:* we get zero if the angle $\Delta/L = \theta$ (Fig. 30-3) corresponds to a 360° phase shift, that is, if Δ is the wavelength λ . This is because the vector contributions form together a complete polygon with zero resultant. The smallest angle that can be resolved by an antenna array of length L is $\theta = \lambda/L$. Notice that the receiving pattern of an antenna such as this is exactly the same as the intensity distribution we would get if we turned the receiver around and made it into a transmitter. This is an example of what is called a *reciprocity principle*. It turns out, in fact, to be generally true for any arrangement of antennas, angles, and so on, that if we first work out what the relative intensities would be in various directions if the receiver were a transmitter instead, then the relative directional sensitivity of a receiver with the same external wiring, the same array of antennas, is the same as the relative intensity of emission would be if it were a transmitter.

Some radio antennas are made in a different way. Instead of having a whole lot of dipoles in a long line, with a lot of feed wires, we may arrange them not in a line but in a curve, and put the receiver at a certain point where it can detect the scattered waves. This curve is cleverly designed so that if the radiowaves are coming down from above, and the wires scatter, making a new wave, the wires are so arranged that the scattered waves reach the receiver all at the same time (Fig. 26-12). In other words, the curve is a *parabola*, and when the source is exactly on its axis, we get a very strong intensity at the focus. In this case we understand very clearly what the resolving power of such an instrument is. The arranging of the antennas on a parabolic curve is not an essential point. It is only a convenient way to get all the signals to the same point with no relative delay and without feed wires. The angle such an instrument can resolve is still $\theta = \lambda/L$, where L is the separation of the first and last antennas. It does not depend on the spacing of the antennas and they may be very close together or in fact be all one piece of metal. Now we are describing a telescope mirror, of course. We have found the resolving power of a telescope! (Sometimes the resolving power is written $\theta = 1.22\lambda/L$, where L is the diameter of the telescope. The reason that it is not exactly λ/L is this: when we worked out that $\theta = \lambda/L$, we assumed that all the lines of dipoles were equal in strength, but when we have a circular telescope, which is the way we usually arrange a telescope, not as much signal comes from the outside edges, because it is not like a square, where we get the same intensity all along a side. We get somewhat less because we are using only part of the telescope there; thus we can appreciate that the effective diameter is a little shorter than the true diameter, and that is what the 1.22 factor tells us. In any case, it seems a little pedantic to put such precision into the resolving power formula.*)

30-5 Colored films; crystals

The above, then, are some of the effects of interference obtained by adding the various waves. But there are a number of other examples, and even though we do not understand the fundamental mechanism yet, we will some day, and we can understand even now how the interference occurs. For example, when a light wave hits a surface of a material with an index n , let us say at normal incidence, some of the light is reflected. The *reason* for the reflection we are not in a position to understand right now; we shall discuss it later. But suppose we know that some of the light is reflected both on entering and leaving a refracting medium. Then, if we look at the reflection of a light source in a thin film, we see the sum of two waves; if the thicknesses are small enough, these two waves will produce an interference, either constructive or destructive, depending on the signs of the phases. It might be, for instance, that for red light, we get an enhanced reflection, but for

* This is because Rayleigh's criterion is a rough idea in the first place. It tells you where it begins to get very hard to tell whether the image was made by one or by two stars. Actually, if sufficiently careful measurements of the exact intensity distribution over the diffracted image spot can be made, the fact that two sources make the spot can be proved even if θ is less than λ/L .

blue light, which has a different wavelength, perhaps we get a destructively interfering reflection, so that we see a bright red reflection. If we change the thickness, i.e., if we look at another place where the film is thicker, it may be reversed, the red interfering and the blue not, so it is bright blue, or green, or yellow, or whatnot. So we see *colors* when we look at thin films and the colors change if we look at different angles, because we can appreciate that the timings are different at different angles. Thus we suddenly appreciate another hundred thousand situations involving the colors that we see on oil films, soap bubbles, etc. at different angles. But the principle is all the same: we are only adding waves at different phases.

As another important application of diffraction, we may mention the following. We used a grating and we saw the diffracted image on the screen. If we had used monochromatic light, it would have been at a certain specific place. Then there were various higher-order images also. From the positions of the images, we could tell how far apart the lines on the grating were, if we knew the wavelength of the light. From the difference in intensity of the various images, we could find out the shape of the grating scratches, whether the grating was made of wires, sawtooth notches, or whatever, *without being able to see them*. This principle is used to discover the positions of the *atoms in a crystal*. The only complication is that a crystal is three-dimensional; it is a repeating three-dimensional array of atoms. We cannot use ordinary light, because we must use something whose wavelength is less than the space between the atoms or we get no effect; so we must use radiation of very short wavelength, i.e., x-rays. So, by shining x-rays into a crystal and by noticing how intense is the reflection in the various orders, we can determine the arrangement of the atoms inside without ever being able to see them with the eye! It is in this way that we know the arrangement of the atoms in various substances, which permitted us to draw those pictures in the first chapter, showing the arrangement of atoms in salt, and so on. We shall later come back to this subject and discuss it in more detail, and therefore we say no more about this most remarkable idea at present.

30-6 Diffraction by opaque screens

Now we come to a very interesting situation. Suppose that we have an opaque sheet with holes in it, and a light on one side of it. We want to know what the intensity is on the other side. What most people say is that the light shines through the holes, and produces an effect on the other side. It will turn out that one gets the right answer, to an excellent approximation, if he assumes that there are sources distributed with uniform density across the open holes, and that the phases of these sources are the same as they would have been if the opaque material were absent. Of course, actually there are *no* sources at the holes, in fact that is the only place that there are *certainly* no sources. Nevertheless, we get the correct diffraction patterns by considering the holes to be the only places that there *are* sources; that is a rather peculiar fact. We shall explain later why this is true, but for now let us just suppose that it is.

In the theory of diffraction there is another kind of diffraction that we should briefly discuss. It is usually not discussed in an elementary course as early as this, only because the mathematical formulas involved in adding these little vectors are a little elaborate. Otherwise it is exactly the same as we have been doing all along. All the interference phenomena are the same; there is nothing very much more advanced involved, only the circumstances are more complicated and it is harder to add the vectors together, that is all.

Suppose that we have light coming in from infinity, casting a shadow of an object. Figure 30-7 shows a screen on which the shadow of an object AB is made by a light source very far away compared with one wavelength. Now we would expect that outside the shadow, the intensity is all bright, and inside it, it is all dark. As a matter of fact, if we plot the intensity as a function of position near the shadow edge, the intensity rises and then overshoots, and wobbles, and oscillates about in a very peculiar manner near this edge (Fig. 30-8). We now shall discuss the reason for this. If we use the theorem that we have not yet proved, then we can

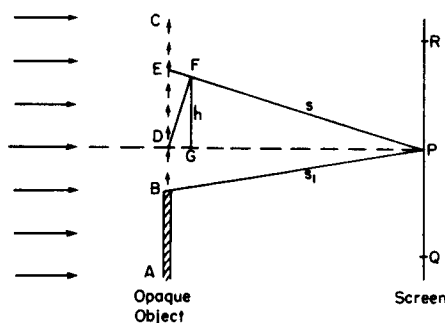


Fig. 30-7. A distant light source casts a shadow of an opaque object on a screen.

replace the actual problem by a set of effective sources uniformly distributed over the open space beyond the object.

We imagine a large number of very closely spaced antennas, and we want the intensity at some point P . That looks just like what we have been doing. Not quite; because our screen is not at infinity. We do not want the intensity at infinity, but at a finite point. To calculate the intensity at some particular place, we have to add the contributions from all the antennas. First there is an antenna at D , exactly opposite P ; if we go up a little bit in angle, let us say a height h , then there is an increase in delay (there is also a change in amplitude because of the change in distance, but this is a very small effect if we are at all far away, and is much less important than the difference in the phases). Now the path difference $EP - DP$ is $h^2/2s$, so that the phase difference is proportional to the *square* of how far we go from D , while in our previous work s was infinite, and the phase difference was *linearly* proportional to h . When the phases are linearly proportional, each vector adds at a constant angle to the next vector. What we now need is a curve which is made by adding a lot of infinitesimal vectors with the requirement that the angle they make shall increase, not linearly, but as the *square* of the length of the curve. To construct that curve involves slightly advanced mathematics, but we can always construct it by actually drawing the arrows and measuring the angles. In any case, we get the marvelous curve (called Cornu's spiral) shown in Fig. 30-8. Now how do we use this curve?

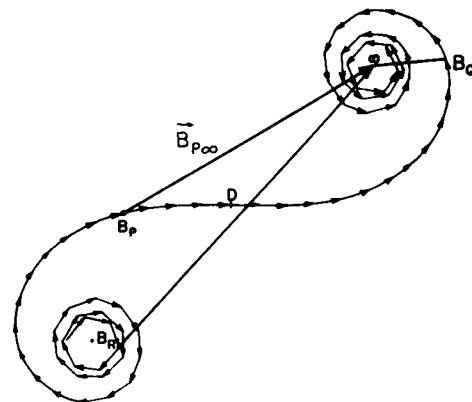


Fig. 30-8. The addition of amplitudes for many in-phase oscillators whose phase delays vary as the square of the distance from point D of the previous figure.

If we want the intensity, let us say, at point P , we add a lot of contributions of different phases from point D on up to infinity, and from D down only to point B_P . So we start at B_P in Fig. 30-8, and draw a series of arrows of ever-increasing angle. Therefore the total contribution above point B_P all goes along the spiraling curve. If we were to stop integrating at some place, then the total amplitude would be a vector from B to that point; in this particular problem we are going to infinity, so the total answer is the vector $B_{P\infty}$. Now the position on the curve which corresponds to point B_P on the object depends upon where point P is located, since point D , the inflection point, always corresponds to the position of point P . Thus, depending upon where P is located above B , the beginning point will fall at various positions on the lower left part of the curve, and the resultant vector $B_{P\infty}$ will have many maxima and minima (Fig. 30-9).

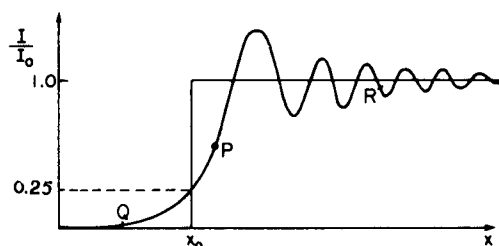


Fig. 30-9. The intensity near the edge of a shadow. The geometrical shadow edge is at x_0 .

On the other hand, if we are at Q , on the other side of P , then we are using only one end of the spiral curve, and not the other end. In other words, we do not even start at D , but at B_Q , so on this side we get an intensity which continuously falls off as Q goes farther into the shadow.

One point that we can immediately calculate with ease, to show that we really understand it, is the intensity exactly opposite the edge. The intensity here is $1/4$ that of the incident light. Reason: Exactly at the edge (so the endpoint B of the arrow is at D in Fig. 30-8) we have half the curve that we would have had if we were far into the bright region. If our point R is far into the light we go from one end of the curve to the other, that is, one full unit vector; but if we are at the edge of the shadow, we have only half the amplitude— $1/4$ the intensity.

In this chapter we have been finding the intensity produced in various directions from various distributions of sources. As a final example we shall derive a formula which we shall need for the next chapter on the theory of the index of refraction. Up to this point relative intensities have been sufficient for our purpose, but this time we shall find the complete formula for the field in the following situation.

30-7 The field of a plane of oscillating charges

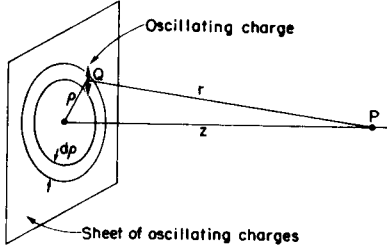


Fig. 30-10. Radiation field of a sheet of oscillating charges.

Suppose that we have a plane full of sources, all oscillating together, with their motion in the plane and all having the same amplitude and phase. What is the field at a finite, but very large, distance away from the plane? (We cannot get very close, of course, because we do not have the right formulas for the field close to the sources.) If we let the plane of the charges be the XY -plane, then we want to find the field at the point P far out on the Z -axis (Fig. 30-10). We suppose that there are η charges per unit area of the plane, and that each one of them has a charge q . All of the charges move with simple harmonic motion, with the same direction, amplitude, and phase. We let the motion of each charge, with respect to its own average position, be $x_0 \cos \omega t$. Or, using the complex notation and remembering that the real part represents the actual motion, the motion can be described by $x_0 e^{i\omega t}$.

Now we find the field at the point P from all of the charges by finding the field there from each charge q , and then adding the contributions from all the charges. We know that the radiation field is proportional to the acceleration of the charge, which is $-\omega^2 x_0 e^{i\omega t}$ (and is the same for every charge). The electric field that we want at the point P due to a charge at the point Q is proportional to the acceleration of the charge q , but we have to remember that the field at the point P at the instant t is given by the acceleration of the charge at the earlier time $t' = t - r/c$, where r/c is the time it takes the waves to travel the distance r from Q to P . Therefore the field at P is proportional to

$$-\omega^2 x_0 e^{i\omega(t-r/c)}. \quad (30.10)$$

Using this value for the acceleration as seen from P in our formula for the electric field at large distances from a radiating charge, we get

$$\left(\text{Electric field at } P \right)_{\text{from charge at } Q} = \frac{q}{4\pi\epsilon_0 c^2} \frac{\omega^2 x_0 e^{i\omega(t-r/c)}}{r} \quad (\text{approx.}). \quad (30.11)$$

Now this formula is not quite right, because we should have used *not* the acceleration of the charge but *its component* perpendicular to the line QP . We shall suppose, however, that the point P is so far away, compared with the distance of the point Q from the axis (the distance ρ in Fig. 30-9), for those changes that we need to take into account, that we can leave out the cosine factor (which would be nearly equal to 1 anyway).

To get the total field at P , we now add the effects of all the charges in the plane. We should, of course, make a *vector* sum. But since the direction of the electric field is nearly the same for all the charges, we may, in keeping with the approximation we have already made, just add the magnitudes of the fields. To our approximation the field at P depends only on the distance r , so all charges at the same r produce equal fields. So we add, first, the fields of those charges in a ring of width $d\rho$ and radius ρ . Then, by taking the integral over all ρ , we will obtain the total field.

The number of charges in the ring is the product of the surface area of the ring, $2\pi\rho d\rho$, and η , the number of charges per unit area. We have, then,

$$\text{Total field at } P = \int \frac{q}{4\pi\epsilon_0 c^2} \frac{\omega^2 x_0 e^{i\omega(t-r/c)}}{r} \cdot \eta \cdot 2\pi\rho d\rho. \quad (30.12)$$

We wish to evaluate this integral from $\rho = 0$ to $\rho = \infty$. The variable t , of course, is to be held fixed while we do the integral, so the only varying quantities are ρ and r . Leaving out all the constant factors, including the factor $e^{i\omega t}$, for the moment, the integral we wish is

$$\int_{\rho=0}^{\rho=\infty} \frac{e^{-i\omega r/c}}{r} \rho d\rho. \quad (30.13)$$

To do this integral we need to use the relation between r and ρ :

$$r^2 = \rho^2 + z^2. \quad (30.14)$$

Since z is independent of ρ , when we take the differential of this equation, we get

$$2r dr = 2\rho d\rho,$$

which is lucky, since in our integral we can replace $\rho d\rho$ by $r dr$ and the r will cancel the one in the denominator. The integral we want is then the simpler one

$$\int_{r=z}^{r=\infty} e^{-i\omega r/c} dr. \quad (30.15)$$

To integrate an exponential is very easy. We divide by the coefficient of r in the exponent and evaluate the exponential at the limits. But the limits of r are not the same as the limits of ρ . When $\rho = 0$, we have $r = z$, so the limits of r are z to infinity. We get for the integral

$$-\frac{c}{i\omega} [e^{-i\infty} - e^{-(i\omega/c)z}], \quad (30.16)$$

where we have written ∞ for $(r/c)\infty$, since they both just mean a very large number!

Now $e^{-i\infty}$ is a mysterious quantity. Its real part, for example, is $\cos(-\infty)$, which, mathematically speaking, is completely indefinite (although we would expect it to be somewhere—or everywhere (?)—between $+1$ and -1 !). But in a *physical* situation, it can mean something quite reasonable, and usually can just be taken to be zero. To see that this is so in our case, we go back to consider again the original integral (30.15).

We can understand (30.15) as a sum of many small complex numbers, each of magnitude Δr , and with the angle $\theta = -\omega r/c$ in the complex plane. We can try to evaluate the sum by a graphical method. In Fig. 30-11 we have drawn the first five pieces of the sum. Each segment of the curve has the length Δr and is placed at the angle $\Delta\theta = -\omega \Delta r/c$ with respect to the preceding piece. The sum for these first five pieces is represented by the arrow from the starting point to the end of the fifth segment. As we continue to add pieces we shall trace out a polygon until we get back to the starting point (approximately) and then start around once more. Adding more pieces, we just go round and round, staying close to a circle whose radius is easily shown to be c/ω . We can see now why the integral does not give a definite answer!

But now we have to go back to the *physics* of the situation. In any real situation the plane of charges *cannot* be infinite in extent, but must sometime stop. If it stopped suddenly, and was exactly circular in shape, our integral would have some value on the circle in Fig. 30-11. If, however, we let the number of charges in the plane gradually taper off at some large distance from the center (or else stop suddenly but in an irregular shape so for larger ρ the entire ring of width $d\rho$ no longer contributes), then the coefficient η in the exact integral would decrease toward zero. Since we are adding smaller pieces but still turning through the same angle, the graph of our integral would then become a curve which is a spiral. The spiral would eventually end up at the center of our original circle, as drawn in Fig. 30-12. The *physically* correct integral is the complex number A in the figure represented by the interval from the starting point to the center of the circle, which is just equal to

$$\frac{c}{i\omega} e^{-i\omega z/c}, \quad (30.17)$$

as you can work out for yourself. This is the same result we would get from Eq. (30.16) if we set $e^{-i\infty} = 0$.

(There is also another reason why the contribution to the integral tapers off for large values of r , and that is the factor we have omitted for the projection of the acceleration on the plane perpendicular to the line PQ .)

We are, of course, interested only in physical situations, so we will take $e^{-i\infty}$ equal to zero. Returning to our original formula (30.12) for the field and putting

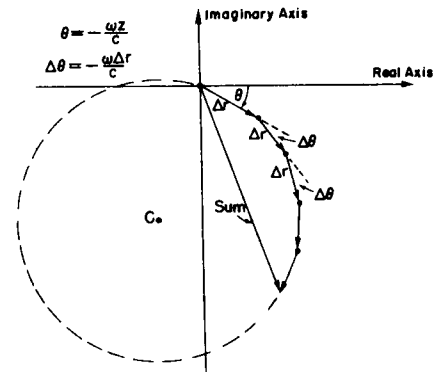


Fig. 30-11. Graphical solution of $\int_z^\infty e^{-i\omega r/c} dr$.

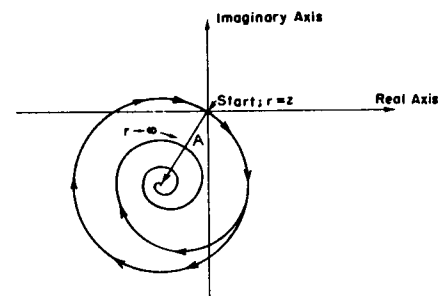


Fig. 30-12. Graphical solution of $\int_z^\infty \eta e^{-i\omega r/c} dr$.

back all of the factors that go with the integral, we have the result

$$\text{Total field at } P = - \frac{\eta q}{2\epsilon_0 c} i\omega x_0 e^{i\omega(t-z/c)} \quad (30.18)$$

(remembering that $1/i = -i$).

It is interesting to note that $(i\omega x_0 e^{i\omega t})$ is just equal to the *velocity* of the charges, so that we can also write the equation for the field as

$$\text{Total field at } P = - \frac{\eta q}{2\epsilon_0 c} [\text{velocity of charges}]_{\text{at } t-z/c}, \quad (30.19)$$

which is a little strange, because the retardation is just by the distance z , which is the shortest distance from P to the plane of charges. But that is the way it comes out—fortunately a rather simple formula. (We may add, by the way, that although our derivation is valid only for distances far from the plane of oscillatory charges, it turns out that the formula (30.18) or (30.19) is correct at any distance z , even for $z < \lambda$.)

The Origin of the Refractive Index

31-1 The index of refraction

We have said before that light goes slower in water than in air, and slower, slightly, in air than in vacuum. This effect is described by the index of refraction n . Now we would like to understand how such a slower velocity could come about. In particular, we should try to see what the relation is to some physical assumptions, or statements, we made earlier, which were the following:

- (a) That the total electric field in any physical circumstance can always be represented by the sum of the fields from all the charges in the universe.
- (b) That the field from a single charge is given by its acceleration evaluated with a retardation at the speed c , *always* (for the *radiation* field).

But, for a piece of glass, you might think: "Oh, no, you should modify all this. You should say it is retarded at the speed c/n ." That, however, is not right, and we have to understand why it is not.

It is approximately true that light or any electrical wave *does appear* to travel at the speed c/n through a material whose index of refraction is n , but the fields are still produced by the motions of *all* the charges—including the charges moving in the material—and with these basic contributions of the field travelling at the ultimate velocity c . Our problem is to understand how the *apparently* slower velocity comes about.

We shall try to understand the effect in a very simple case. A source which we shall call "the *external source*" is placed a large distance away from a thin plate of transparent material, say glass. We inquire about the field at a large distance on the opposite side of the plate. The situation is illustrated by the diagram of Fig. 31-1, where S and P are imagined to be very far away from the plate. According to the principles we have stated earlier, an electric field anywhere that is far from all moving charges is the (vector) sum of the fields produced by the external source (at S) and the fields produced by *each* of the charges in the plate of glass, *every one with its proper retardation at the velocity c* . Remember that the contribution of each charge is not changed by the presence of the other charges. These are our basic principles. The field at P can be written thus:

$$\mathbf{E} = \sum_{\text{all charges}} \mathbf{E}_{\text{each charge}} \quad (31.1)$$

or

$$\mathbf{E} = \mathbf{E}_s + \sum_{\text{all other charges}} \mathbf{E}_{\text{each charge}}, \quad (31.2)$$

where \mathbf{E}_s is the field due to the source alone and would be precisely the field at P if there were no material present. We expect the field at P to be different from \mathbf{E}_s if there are any other moving charges.

Why should there be charges moving in the glass? We know that all material consists of atoms which contain electrons. When the electric field of the source acts on these atoms it drives the electrons up and down, because it exerts a force on the electrons. And moving electrons generate a field—they constitute new radiators. These new radiators are related to the source S , because they are driven by the field of the source. The total field is not just the field of the source S , but it is modified by the additional contribution from the other moving charges. This means that the field is not the same as the one which was there before the glass was there, but is modified, and it turns out that it is modified in such a way that

31-1 The index of refraction

31-2 The field due to the material

31-3 Dispersion

31-4 Absorption

31-5 The energy carried by an electric wave

31-6 Diffraction of light by a screen

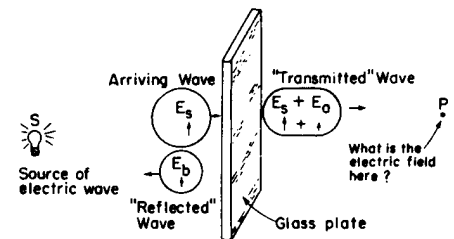


Fig. 31-1. Electric waves passing through a layer of transparent material.

the field inside the glass appears to be moving at a different speed. That is the idea which we would like to work out quantitatively.

Now this is, in the exact case, pretty complicated, because although we have said that all the other moving charges are driven by the source field, that is not quite true. If we think of a particular charge, it feels not only the source, but like anything else in the world, it feels *all* of the charges that are moving. It feels, in particular, the charges that are moving somewhere else in the glass. So the total field which is acting on a *particular charge* is a combination of the fields from the other charges, *whose motions depend on what this particular charge is doing!* You can see that it would take a complicated set of equations to get the complete and exact formula. It is so complicated that we postpone this problem until next year.

Instead we shall work out a very simple case in order to understand all the physical principles very clearly. We take a circumstance in which the effects from the other atoms is very small relative to the effects from the source. In other words, we take a material in which the total field is not modified very much by the motion of the other charges. That corresponds to a material in which the index of refraction is very close to 1, which will happen, for example, if the density of the atoms is very low. Our calculation will be valid for any case in which the index is for any reason very close to 1. In this way we shall avoid the complications of the most general, complete solution.

Incidentally, you should notice that there is another effect caused by the motion of the charges in the plate. These charges will also radiate waves back toward the source S . This backward-going field is the light we see reflected from the surfaces of transparent materials. It does not come from just the surface. The backward radiation comes from everywhere in the interior, but it turns out that the total effect is equivalent to a reflection from the surfaces. These reflection effects are beyond our approximation at the moment because we shall be limited to a calculation for a material with an index so close to 1 that very little light is reflected.

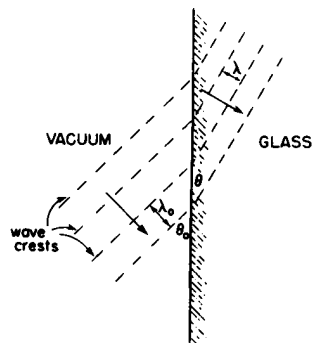


Fig. 31-2. Relation between refraction and velocity change.

Before we proceed with our study of how the index of refraction comes about, we should understand that all that is required to understand refraction is to understand why the apparent wave *velocity* is different in different materials. The *bending* of light rays comes about just *because* the effective speed of the waves is different in the materials. To remind you how that comes about we have drawn in Fig. 31-2 several successive crests of an electric wave which arrives from a vacuum onto the surface of a block of glass. The arrow perpendicular to the wave crests indicates the direction of travel of the wave. Now all oscillations in the wave must have the same *frequency*. (We have seen that driven oscillations have the same frequency as the driving source.) This means, also, that the wave crests for the waves on both sides of the surface must have the *same spacing along the surface* because they must travel together, so that a charge sitting at the boundary will feel only one frequency. The *shortest* distance between crests of the wave, however, is the wavelength which is the velocity divided by the frequency. On the vacuum side it is $\lambda_0 = 2\pi c/\omega$, and on the other side it is $\lambda = 2\pi v/\omega$ or $2\pi c/\omega n$, if $v = c/n$ is the velocity of the wave. From the figure we can see that the only way for the waves to “fit” properly at the boundary is for the waves in the material to be travelling at a different angle with respect to the surface. From the geometry of the figure you can see that for a “fit” we must have $\lambda_0/\sin \theta_0 = \lambda/\sin \theta$, or $\sin \theta_0/\sin \theta = n$, which is Snell’s law. We shall, for the rest of our discussion, consider only why light has an effective speed of c/n in material of index n , and no longer worry, in this chapter, about the bending of the light direction.

We go back now to the situation shown in Fig. 31-1. We see that what we have to do is to calculate the field produced at P by all the oscillating charges in the glass plate. We shall call this part of the field E_a , and it is just the sum written as the second term in Eq. (31.2). When we add it to the term E_s , due to the source, we will have the total field at P .

This is probably the most complicated thing that we are going to do this year, but it is complicated only in that there are many pieces that have to be put together; each piece, however, is very simple. Unlike other derivations where we say, "Forget the derivation, just look at the answer!" in this case we do not need the answer so much as the derivation. In other words, the thing to understand now is the physical machinery for the production of the index.

To see where we are going, let us first find out what the "correction field" E_a would have to be if the total field at P is going to look like radiation from the source that is slowed down while passing through the thin plate. If the plate had no effect on it, the field of a wave travelling to the right (along the z -axis) would be

$$E_s = E_0 \cos \omega(t - z/c) \quad (31.3)$$

or, using the exponential notation,

$$E_s = E_0 e^{i\omega(t-z/c)}. \quad (31.4)$$

Now what would happen if the wave travelled more slowly in going through the plate? Let us call the thickness of the plate Δz . If the plate were not there the wave would travel the distance Δz in the time $\Delta z/c$. But if it appears to travel at the speed c/n then it should take the longer time $n \Delta z/c$ or the *additional* time $\Delta t = (n - 1) \Delta z/c$. After that it would continue to travel at the speed c again. We can take into account the extra delay in getting through the plate by replacing t in Eq. (31.4) by $(t - \Delta t)$ or by $[t - (n - 1) \Delta z/c]$. So the wave after insertion of the plate should be written

$$E_{\text{after plate}} = E_0 e^{i\omega[t - (n-1)\Delta z/c - z/c]}. \quad (31.5)$$

We can also write this equation as

$$E_{\text{after plate}} = e^{-i\omega(n-1)\Delta z/c} E_0 e^{i\omega(t-z/c)}, \quad (31.6)$$

which says that the wave after the plate is obtained from the wave which could exist without the plate, i. e., from E_s , by multiplying by the factor $e^{-i\omega(n-1)\Delta z/c}$. Now we know that multiplying an oscillating function like $e^{i\omega t}$ by a factor $e^{i\theta}$ just says that we change the phase of the oscillation by the angle θ , which is, of course, what the extra delay in passing through the thickness Δz has done. It has retarded the phase by the amount $\omega(n - 1) \Delta z/c$ (retarded, because of the minus sign in the exponent).

We have said earlier that the plate should *add* a field E_a to the original field $E_s = E_0 e^{i\omega(t-z/c)}$, but we have found instead that the effect of the plate is to *multiply* the field by a factor which shifts its phase. However, that is really all right because we can get the same result by adding a suitable complex number. It is particularly easy to find the right number to add in the case that Δz is small, for you will remember that if x is a small number then e^x is nearly equal to $(1 + x)$. We can write, therefore,

$$e^{-i\omega(n-1)\Delta z/c} = 1 - i\omega(n - 1) \Delta z/c. \quad (31.7)$$

Using this equality in Eq. (31.6), we have

$$E_{\text{after plate}} = \underbrace{E_0 e^{i\omega(t-z/c)}}_{E_s} - \underbrace{\frac{i\omega(n-1)\Delta z}{c} E_0 e^{i\omega(t-z/c)}}_{E_a}. \quad (31.8)$$

The first term is just the field from the source, and the second term must just be equal to E_a , the field produced to the right of the plate by the oscillating charges of the plate—expressed here in terms of the index of refraction n , and depending, of course, on the strength of the wave from the source.

What we have been doing is easily visualized if we look at the complex number diagram in Fig. 31-3. We first draw the number E_s (we chose some values for z and t so that E_s comes out horizontal, but this is not necessary). The delay due to

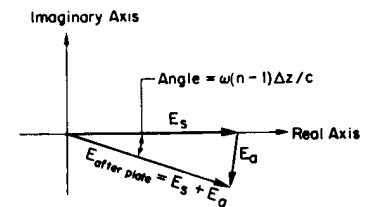


Fig. 31-3. Diagram for the transmitted wave at a particular t and z .

slowing down in the plate would delay the phase of this number, that is, it would rotate E_s through a negative angle. But this is equivalent to adding the small vector E_a at roughly right angles to E_s . But that is just what the factor $-i$ means in the second term of Eq. (31.8). It says that if E_s is real, then E_a is negative imaginary or that, in general, E_s and E_a make a right angle.

31-2 The field due to the material

We now have to ask: Is the field E_a obtained in the second term of Eq. (31.8) the kind we would expect from oscillating charges in the plate? If we can show that it is, we will then have calculated what the index n should be! [Since n is the only nonfundamental number in Eq. (31.8).] We turn now to calculating what field E_a the charges in the material will produce. (To help you keep track of the many symbols we have used up to now, and will be using in the rest of our calculation, we have put them all together in Table 31-1.)

Table 31-1

Symbols used in the calculations

E_s	= field from the source
E_a	= field produced by charges in the plate
Δz	= thickness of the plate
z	= perpendicular distance from the plate
n	= index of refraction
ω	= frequency (angular) of the radiation
N	= number of charges per unit volume in the plate
η	= number of charges per unit area of the plate
q_e	= charge on an electron
m	= mass of an electron
ω_0	= resonant frequency of an electron bound in an atom

If the source S (of Fig. 31-1) is far off to the left, then the field E_s will have the same phase everywhere on the plate, so we can write that in the neighborhood of the plate

$$E_s = E_0 e^{i(\omega t - z/c)} \quad (31.9)$$

Right at the plate, where $z = 0$, we will have

$$E_s = E_0 e^{i\omega t} \text{ (at the plate)} \quad (31.10)$$

Each of the electrons in the atoms of the plate will feel this electric field and will be driven up and down (we assume the direction of E_0 is vertical) by the electric force qE . To find what motion we expect for the electrons, we will assume that the atoms are little oscillators, that is, that the electrons are fastened elastically to the atoms, which means that if a force is applied to an electron its displacement from its normal position will be proportional to the force.

You may think that this is a funny model of an atom if you have heard about electrons whirling around in orbits. But that is just an oversimplified picture. The correct picture of an atom, which is given by the theory of wave mechanics, says that, *so far as problems involving light are concerned*, the electrons behave as though they were held by springs. So we shall suppose that the electrons have a linear restoring force which, together with their mass m , makes them behave like little oscillators, with a resonant frequency ω_0 . We have already studied such oscillators, and we know that the equation of their motion is written this way:

$$m \left(\frac{d^2 x}{dt^2} + \omega_0^2 x \right) = F, \quad (31.11)$$

where F is the driving force.

For our problem, the driving force comes from the electric field of the wave from the source, so we should use

$$F = q_e E_s = q_e E_0 e^{i\omega t}, \quad (31.12)$$

where q_e is the electric charge on the electron and for E_s we use the expression $E_s = E_0 e^{i\omega t}$ from (31.10). Our equation of motion for the electron is then

$$m \left(\frac{d^2 x}{dt^2} + \omega_0^2 x \right) = q_e E_0 e^{i\omega t}. \quad (31.13)$$

We have solved this equation before, and we know that the solution is

$$x = x_0 e^{i\omega t}, \quad (31.14)$$

where, by substituting in (31.13), we find that

$$x_0 = \frac{q_e E_0}{m(\omega_0^2 - \omega^2)}, \quad (31.15)$$

so that

$$x = \frac{q_e E_0}{m(\omega_0^2 - \omega^2)} e^{i\omega t}. \quad (31.16)$$

We have what we needed to know—the motion of the electrons in the plate. And it is the same for every electron, except that the mean position (the “zero” of the motion) is, of course, different for each electron.

Now we are ready to find the field E_a that these atoms produce at the point P , because we have already worked out (at the end of Chapter 30) what field is produced by a sheet of charges that all move together. Referring back to Eq. (30.19), we see that the field E_a at P is just a negative constant times the velocity of the charges retarded in time the amount z/c . Differentiating x in Eq. (31.16) to get the velocity, and sticking in the retardation [or just putting x_0 from (31.15) into (30.18)] yields

$$E_a = - \frac{\eta q_e}{2\epsilon_0 c} \left[i\omega \frac{q_e E_0}{m(\omega_0^2 - \omega^2)} e^{i\omega(t-z/c)} \right]. \quad (31.17)$$

Just as we expected, the driven motion of the electrons produced an extra wave which travels to the right (that is what the factor $e^{i\omega(t-z/c)}$ says), and the amplitude of this wave is proportional to the number of atoms per unit area in the plate (the factor η) and also proportional to the strength of the source field (the factor E_0). Then there are some factors which depend on the atomic properties (q_e , m , and ω_0), as we should expect.

The most important thing, however, is that this formula (31.17) for E_a looks very much like the expression for E_a that we got in Eq. (31.8) by saying that the original wave was delayed in passing through a material with an index of refraction n . The two expressions will, in fact, be identical if

$$(n - 1) \Delta z = \frac{\eta q_e^2}{2\epsilon_0 m(\omega_0^2 - \omega^2)}. \quad (31.18)$$

Notice that both sides are proportional to Δz , since η , which is the number of atoms *per unit area*, is equal to $N \Delta z$, where N is the number of atoms *per unit volume* of the plate. Substituting $N \Delta z$ for η and cancelling the Δz , we get our main result, a formula for the index of refraction in terms of the properties of the atoms of the material—and of the frequency of the light:

$$n = 1 + \frac{N q_e^2}{2\epsilon_0 m(\omega_0^2 - \omega^2)}. \quad (31.19)$$

This equation gives the “explanation” of the index of refraction that we wished to obtain.

31-3 Dispersion

Notice that in the above process we have obtained something very interesting. For we have not only a number for the index of refraction which can be computed from the basic atomic quantities, but we have also learned how the index of refraction should vary with the frequency ω of the light. This is something we would never understand from the simple statement that "light travels slower in a transparent material." We still have the problem, of course, of knowing how many atoms per unit volume there are, and what is their natural frequency ω_0 . We do not know this just yet, because it is different for every different material, and we cannot get a general theory of that now. Formulation of a general theory of the properties of different substances—their natural frequencies, and so on—is possible only with quantum atomic mechanics. Also, different materials have different properties and different indexes, so we cannot expect, anyway, to get a general formula for the index which will apply to all substances.

However, we shall discuss the formula we have obtained, in various possible circumstances. First of all, for most ordinary gases (for instance, for air, most colorless gases, hydrogen, helium, and so on) the natural frequencies of the electron oscillators correspond to ultraviolet light. These frequencies are higher than the frequencies of visible light, that is, ω_0 is much larger than ω of visible light, and to a first approximation, we can disregard ω^2 in comparison with ω_0^2 . Then we find that the index is nearly constant. So for a gas, the index is nearly constant. This is also true for most other transparent substances, like glass. If we look at our expression a little more closely, however, we notice that as ω rises, taking a little bit more away from the denominator, the index also rises. So n rises slowly with frequency. The index is higher for blue light than for red light. That is the reason why a prism bends the light more in the blue than in the red.

The phenomenon that the index depends upon the frequency is called the phenomenon of *dispersion*, because it is the basis of the fact that light is "dispersed" by a prism into a spectrum. The equation for the index of refraction as a function of frequency is called a *dispersion equation*. So we have obtained a dispersion equation. (In the past few years "dispersion equations" have been finding a new use in the theory of elementary particles.)

Our dispersion equation suggests other interesting effects. If we have a natural frequency ω_0 which lies in the visible region, or if we measure the index of refraction of a material like glass in the ultraviolet, where ω gets near ω_0 , we see that at frequencies very close to the natural frequency the index can get enormously large, because the denominator can go to zero. Next, suppose that ω is greater than ω_0 . This would occur, for example, if we take a material like glass, say, and shine x-ray radiation on it. In fact, since many materials which are opaque to visible light, like graphite for instance, are transparent to x-rays, we can also talk about the index of refraction of carbon for x-rays. All the natural frequencies of the carbon atoms would be much lower than the frequency we are using in the x-rays, since x-ray radiation has a very high frequency. The index of refraction is that given by our dispersion equation if we set ω_0 equal to zero (we neglect ω_0^2 in comparison with ω^2).

A similar situation would occur if we beam radiowaves (or light) on a gas of free electrons. In the upper atmosphere electrons are liberated from their atoms by ultraviolet light from the sun and they sit up there as free electrons. For free electrons $\omega_0 = 0$ (there is no elastic restoring force). Setting $\omega_0 = 0$ in our dispersion equation yields the correct formula for the index of refraction for radiowaves in the stratosphere, where N is now to represent the density of free electrons (number per unit volume) in the stratosphere. But let us look again at the equation, if we beam x-rays on matter, or radiowaves (or any electric waves) on free electrons the term $(\omega_0^2 - \omega^2)$ becomes *negative*, and we obtain the result that n is *less than one*. That means that the effective speed of the waves in the substance is *faster* than c ! Can that be correct?

It is correct. In spite of the fact that it is said that you cannot send signals any faster than the speed of light, it is nevertheless true that the index of refraction of materials at a particular frequency can be either greater or less than 1. This

just means that the *phase shift* which is produced by the scattered light can be either positive or negative. It can be shown, however, that the speed at which you can send a *signal* is not determined by the index at one frequency, but depends on what the index is at *many* frequencies. What the index tells us is the speed at which the *nodes* (or crests) of the wave travel. The *node* of a wave is not a signal by itself. In a perfect wave, which has no modulations of any kind, i.e., which is a steady oscillation, you cannot really say when it “starts,” so you cannot use it for a timing signal. In order to send a *signal* you have to change the wave somehow, make a notch in it, make it a little bit fatter or thinner. That means that you have to have more than one frequency in the wave, and it can be shown that the speed at which *signals* travel is not dependent upon the index alone, but upon the way that the index changes with the frequency. This subject we must also delay (until Chapter 48). Then we will calculate for you the actual speed of *signals* through such a piece of glass, and you will see that it will not be faster than the speed of light, although the nodes, which are mathematical points, do travel faster than the speed of light.

Just to give a slight hint as to how that happens, you will note that the real difficulty has to do with the fact that the responses of the charges are opposite to the field, i.e., the sign has gotten reversed. Thus in our expression for x (Eq. 31.16) the displacement of the charge is in the direction opposite to the driving field, because $(\omega_0^2 - \omega^2)$ is negative for small ω_0 . The formula says that when the electric field is pulling in one direction, the charge is moving in the opposite direction.

How does the charge happen to be going in the opposite direction? It certainly does not start off in the opposite direction when the field is first turned on. When the motion first starts there is a transient, which settles down after awhile, and only *then* is the phase of the oscillation of the charge opposite to the driving field. And it is then that the *phase* of the transmitted field can appear to be *advanced* with respect to the source wave. It is this *advance in phase* which is meant when we say that the “phase velocity” or velocity of the nodes is greater than c . In Fig. 31-4 we give a schematic idea of how the waves might look for a case where the wave is suddenly turned on (to make a signal). You will see from the diagram that the *signal* (i.e., the *start* of the wave) is *not* earlier for the wave which ends up with an advance in phase.

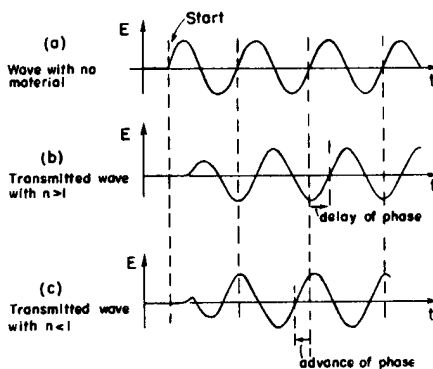


Fig. 31-4. Wave “signals.”

Let us now look again at our dispersion equation. We should remark that our analysis of the refractive index gives a result that is somewhat simpler than you would actually find in nature. To be completely accurate we must add some refinements. First, we should expect that our model of the atomic oscillator should have some damping force (otherwise once started it would oscillate forever, and we do not expect that to happen). We have worked out before (Eq. 23.8) the motion of a damped oscillator and the result is that the denominator in Eq. (31.16), and therefore in (31.19), is changed from $(\omega_0^2 - \omega^2)$ to $(\omega_0^2 - \omega^2 + i\gamma\omega)$, where γ is the damping coefficient.

We need a second modification to take into account the fact that there are several resonant frequencies for a particular kind of atom. It is easy to fix up our

dispersion equation by imagining that there are several different kinds of oscillators, but that each oscillator acts separately, and so we simply add the contributions of all the oscillators. Let us say that there are N_k electrons per unit of volume, whose natural frequency is ω_k and whose damping factor is γ_k . We would then have for our dispersion equation

$$n = 1 + \frac{q_e^2}{2\epsilon_0 m} \sum_k \frac{N_k}{\omega_k^2 - \omega^2 + i\gamma_k \omega}. \quad (31.20)$$

We have, finally, a complete expression which describes the index of refraction that is observed for many substances.* The index described by this formula varies with frequency roughly like the curve shown in Fig. 31-5.

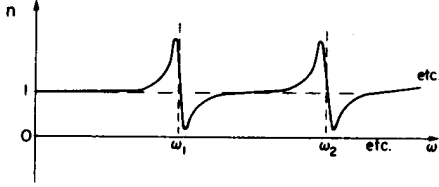


Fig. 31-5. The index of refraction as a function of frequency.

You will note that so long as ω is not too close to one of the resonant frequencies, the slope of the curve is positive. Such a positive slope is called “normal” dispersion (because it is clearly the most common occurrence). Very near the resonant frequencies, however, there is a small range of ω 's for which the slope is negative. Such a negative slope is often referred to as “anomalous” (meaning abnormal) dispersion, because it seemed unusual when it was first observed, long before anyone even knew there were such things as electrons. From our point of view both slopes are quite “normal”!

31-4 Absorption

Perhaps you have noticed something a little strange about the last form (Eq. 31.20) we obtained for our dispersion equation. Because of the term $i\gamma$ we put in to take account of damping, the index of refraction is now a *complex number*! What does *that* mean? By working out what the real and imaginary parts of n are we could write

$$n = n' - in'', \quad (31.21)$$

where n' and n'' are real numbers. (We use the minus sign in front of the in'' because then n'' will turn out to be a positive number, as you can show for yourself.)

We can see what such a complex index means by going back to Eq. (31.6), which is the equation of the wave after it goes through a plate of material with an index n . If we put our complex n into this equation, and do some rearranging, we get

$$E_{\text{after plate}} = \underbrace{e^{-\omega n'' \Delta z / c}}_A \underbrace{e^{-i\omega(n' - 1)\Delta z / c} E_0 e^{i\omega(t - z/c)}}_B \quad (31.22)$$

The last factors, marked B in Eq. (31.22), are just the form we had before, and again describe a wave whose phase has been delayed by the angle $\omega(n' - 1)\Delta z/c$ in traversing the material. The first term (A) is new and is an exponential factor with a *real* exponent, because there were two i 's that cancelled. Also, the exponent is negative, so the factor is a real number less than one. It describes a *decrease* in the magnitude of the field and, as we should expect, by an amount which is more the larger Δz is. As the wave goes through the material, it is weakened. The material is “absorbing” part of the wave. The wave comes out the other side with less energy. We should not be surprised at this, because the damping we put in for the oscillators is indeed a friction force and must be expected to cause a loss of energy. We see that the imaginary part n'' of a complex index of refraction represents an absorption (or “attenuation”) of the wave. In fact, n'' is sometimes referred to as the “absorption index.”

We may also point out that an imaginary part to the index n corresponds to bending the arrow E_a in Fig. 31-3 toward the origin. It is clear why the transmitted field is then decreased.

* Actually, although in quantum mechanics Eq. (31.20) is still valid, its interpretation is somewhat different. In quantum mechanics even an atom with one electron, like hydrogen, has several resonant frequencies. Therefore N_k is not really the number of electrons having the frequency ω_k , but is replaced instead by Nf_k , where N is the number of atoms per unit volume and f_k (called the oscillator strength) is a factor that tells how strongly the atom exhibits each of its resonant frequencies ω_k .

Normally, for instance as in glass, the absorption of light is very small. This is to be expected from our Eq. (31.20), because the imaginary part of the denominator, $i\gamma_k\omega$, is much smaller than the term $(\omega_k^2 - \omega^2)$. But if the light frequency ω is very close to ω_k then the resonance term $(\omega_k^2 - \omega^2)$ can become small compared with $i\gamma_k\omega$ and the index becomes almost completely imaginary. The absorption of the light becomes the dominant effect. It is just this effect that gives the dark lines in the spectrum of light which we receive from the sun. The light from the solar surface has passed through the sun's atmosphere (as well as the earth's), and the light has been strongly absorbed at the resonant frequencies of the atoms in the solar atmosphere.

The observation of such spectral lines in the sunlight allows us to tell the resonant frequencies of the atoms and hence the chemical composition of the sun's atmosphere. The same kind of observations tell us about the materials in the stars. From such measurements we know that the chemical elements in the sun and in the stars are the same as those we find on the earth.

31-5 The energy carried by an electric wave

We have seen that the imaginary part of the index means absorption. We shall now use this knowledge to find out how much energy is carried by a light wave. We have given earlier an argument that the energy carried by light is proportional to $\overline{E^2}$, the time average of the square of the electric field in the wave. The decrease in E due to absorption must mean a loss of energy, which would go into some friction of the electrons and, we might guess, would end up as heat in the material.

If we consider the light arriving on a unit area, say one square centimeter, of our plate in Fig. 31-1, then we can write the following energy equation (if we assume that energy is conserved, as we *do!*):

$$\text{Energy in per sec} = \text{energy out per sec} + \text{work done per sec.} \quad (31.23)$$

For the first term we can write $\alpha\overline{E_s^2}$, where α is the as yet unknown constant of proportionality which relates the average value of E^2 to the energy being carried. For the second term we must include the part from the radiating atoms of the material, so we should use $\alpha(\overline{E_s + E_a})^2$, or (evaluating the square) $\alpha(\overline{E_s^2} + 2\overline{E_s E_a} + \overline{E_a^2})$.

All of our calculations have been made for a thin layer of material whose index is not too far from 1, so that E_a would always be much less than E_s (just to make the calculations easier). In keeping with our approximations, we should, therefore, leave out the term $\overline{E_a^2}$, because it is much smaller than $\overline{E_s E_a}$. You may say: "Then you should leave out $\overline{E_s E_a}$ also, because it is much smaller than $\overline{E_s^2}$." It is true that $\overline{E_s E_a}$ is much smaller than $\overline{E_s^2}$, but we must keep $\overline{E_s E_a}$ or our approximation will be the one that would apply if we neglected the presence of the material completely! One way of checking that our calculations are consistent is to see that we always keep terms which are proportional to $N\Delta z$, the area density of atoms in the material, but we leave out terms which are proportional to $(N\Delta z)^2$ or any higher power of $N\Delta z$. Ours is what should be called a "low-density approximation."

In the same spirit, we might remark that our energy equation has neglected the energy in the reflected wave. But that is OK because this term, too, is proportional to $(N\Delta z)^2$, since the amplitude of the reflected wave is proportional to $N\Delta z$.

For the last term in Eq. (31.23) we wish to compute the rate at which the incoming wave is doing work on the electrons. We know that work is force times distance, so the *rate* of doing work (also called power) is the force times the velocity. It is really $\mathbf{F} \cdot \mathbf{V}$, but we do not need to worry about the dot product when the velocity and force are along the same direction as they are here (except for a possible minus sign). So for each atom we take $q_e \overline{E_s} v$ for the average rate of

doing work. Since there are $N \Delta z$ atoms in a unit area, the last term in Eq. (31.23) should be $N \Delta z q_e \overline{E_s v}$. Our energy equation now looks like

$$\alpha \overline{E_s^2} = \alpha \overline{E_s^2} + 2\alpha \overline{E_s E_a} + N \Delta z q_e \overline{E_s v}. \quad (31.24)$$

The $\overline{E_s^2}$ terms cancel, and we have

$$2\alpha \overline{E_s E_a} = N \Delta z q_e \overline{E_s v}. \quad (31.25)$$

We now go back to Eq. (30.19), which tells us that for large z

$$E_a = \frac{N \Delta z q_e}{2\epsilon_0 c} v(\text{ret by } z/c) \quad (31.26)$$

(recalling that $\eta = N \Delta z$). Putting Eq. (31.26) into the left-hand side of (31.25), we get

$$2\alpha \frac{N \Delta z q_e}{2\epsilon_0 c} \overline{E_s(\text{at } z) \cdot v(\text{ret by } z/c)}.$$

However, $E_s(\text{at } z)$ is E_s (at atoms) retarded by z/c . Since the average is independent of time, it is the same now as retarded by z/c , or is $\overline{E_s}$ (at atom) $\cdot v$, the same average that appears on the right-hand side of (31.25). The two sides are therefore equal if

$$\frac{\alpha}{\epsilon_0 c} = 1, \quad \text{or} \quad \alpha = \epsilon_0 c. \quad (31.27)$$

We have discovered that if energy is to be conserved, the energy carried in an electric wave per unit area and per unit time (or what we have called the *intensity*) must be given by $\epsilon_0 c \overline{E^2}$. If we call the intensity \overline{S} , we have

$$\overline{S} = \left\{ \begin{array}{c} \text{intensity} \\ \text{or} \\ \text{energy/area/time} \end{array} \right\} = \epsilon_0 c \overline{E^2}, \quad (31.28)$$

where the *bar* means the *time average*. We have a nice bonus result from our theory of the refractive index!

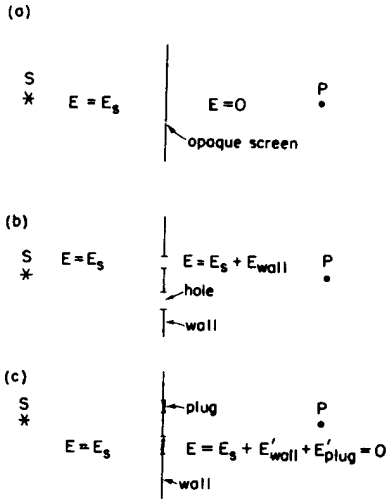


Fig. 31-6. Diffraction by a screen.

31-6 Diffraction of light by a screen

It is now a good time to take up a somewhat different matter which we can handle with the machinery of this chapter. In the last chapter we said that when you have an opaque screen and the light can come through some holes, the distribution of intensity—the diffraction pattern—could be obtained by imagining instead that the holes are replaced by sources (oscillators) uniformly distributed over the hole. In other words, the diffracted wave is the same as though the hole were a new source. We have to explain the reason for that, because the hole is, of course, just where there are *no* sources, where there are *no* accelerating charges.

Let us first ask: “What is an opaque screen?” Suppose we have a completely opaque screen between a source S and an observer at P , as in Fig. 31-6(a). If the screen is “opaque” there is no field at P . Why is there no field there? According to the basic principles we should obtain the field at P as the field E_s of the source delayed, plus the field from all the other charges around. But, as we have seen above, the charges in the screen will be set in motion by the field E_s , and these motions generate a new field which, if the screen is opaque, must *exactly* cancel the field E_s on the back side of the screen. You say: “What a miracle that it balances *exactly*! Suppose it was not exactly right!” If it were not exactly right (remember that this opaque screen has some thickness), the field toward the rear part of the screen would not be exactly zero. So, not being zero, it would set into motion some other charges in the material of the screen, and thus make a little more field, trying to get the total balanced out. So if we make the screen thick enough, there is no residual field, because there is enough opportunity to finally get the thing quieted down. In terms of our formulas above we would say that the

screen has a large and imaginary index, so the wave is absorbed exponentially as it goes through. You know, of course, that a thin enough sheet of the most opaque material, even gold, is transparent.

Now let us see what happens with an opaque screen which has holes in it, as in Fig. 31-6(b). What do we expect for the field at P ? The field at P can be represented as a sum of two parts—the field due to the source S plus the field due to the wall, i.e., due to the motions of the charges in the walls. We might expect the motions of the charges in the walls to be complicated, but we can find out *what fields they produce* in a rather simple way.

Suppose that we were to take the same screen, but plug up the holes, as indicated in part (c) of the figure. We imagine that the plugs are of exactly the same material as the wall. Mind you, the plugs go where the holes were in case (b). Now let us calculate the field at P . The field at P is certainly zero in case (c), but it is *also* equal to the field from the source plus the field due to all the motions of the atoms in the walls and in the plugs. We can write the following equations:

$$\text{Case (b): } E_{\text{at } P} = E_s + E_{\text{wall}},$$

$$\text{Case (c): } E'_{\text{at } P} = 0 = E_s + E'_{\text{wall}} + E'_{\text{plug}},$$

where the primes refer to the case where the plugs are in place, but E_s is, of course, the same in both cases. Now if we subtract the two equations, we get

$$E_{\text{at } P} = (E_{\text{wall}} - E'_{\text{wall}}) - E'_{\text{plug}}.$$

Now if the holes are not too small (say many wavelengths across), we would not expect the presence of the plugs to change the fields which arrive at the walls except possibly for a little bit around the edges of the holes. Neglecting this small effect, we can set $E_{\text{wall}} = E'_{\text{wall}}$ and obtain that

$$E_{\text{at } P} = -E'_{\text{plug}}.$$

We have the result that the field at P *when there are holes* in a screen (case b) is the same (except for sign) as the field that is produced by *that part* of a complete opaque wall which is *located where the holes are!* (The sign is not too interesting, since we are usually interested in intensity which is proportional to the square of the field.) It seems like an amazing backwards-forwards argument. It is, however, not only true (approximately for not too small holes), but useful, and is the justification for the usual theory of diffraction.

The field E'_{plug} is computed in any particular case by remembering that the motion of the charges *everywhere* in the screen is just that which will cancel out the field E_s on the back of the screen. Once we know these motions, we add the radiation fields at P due just to the charges in the plugs.

We remark again that this theory of diffraction is only approximate, and will be good only if the holes are not too small. For holes which are too small the E'_{plug} term will be small and then the difference between E'_{wall} and E_{wall} (which difference we have taken to be zero) may be comparable to or larger than the small E'_{plug} term, and our approximation will no longer be valid.

Radiation Damping. Light Scattering

32-1 Radiation resistance

In the last chapter we learned that when a system is oscillating, energy is carried away, and we deduced a formula for the energy which is radiated by an oscillating system. If we know the electric field, then the average of the square of the field times $\epsilon_0 c$ is the amount of energy that passes per square meter per second through a surface normal to the direction in which the radiation is going:

$$S = \epsilon_0 c \langle E^2 \rangle. \quad (32.1)$$

Any oscillating charge radiates energy; for instance, a driven antenna radiates energy. If the system radiates energy, then in order to account for the conservation of energy we must find that power is being delivered along the wires which lead into the antenna. That is, to the driving circuit the antenna acts like a *resistance*, or a place where energy can be “lost” (the energy is not really lost, it is really radiated out, but so far as the circuit is concerned, the energy is lost). In an ordinary resistance, the energy which is “lost” passes into heat; in this case the energy which is “lost” goes out into space. But from the standpoint of circuit theory, without considering *where* the energy goes, the net effect on the circuit is the same—energy is “lost” from that circuit. Therefore the antenna appears to the generator as having a resistance, even though it may be made with perfectly good copper. In fact, if it is well built it will appear as almost a pure resistance, with very little inductance or capacitance, because we would like to radiate as much energy as possible out of the antenna. This resistance that an antenna shows is called the *radiation resistance*.

If a current I is going to the antenna, then the average rate at which power is delivered to the antenna is the average of the square of the current times the resistance. The rate at which power is *radiated* by the antenna is proportional to the square of the current in the antenna, of course, because all the fields are proportional to the currents, and the energy liberated is proportional to the square of the field. The coefficient of proportionality between radiated power and $\langle I^2 \rangle$ is the radiation resistance.

An interesting question is, what is this radiation resistance due to? Let us take a simple example: let us say that currents are driven up and down in an antenna. We find that we have to put work in, if the antenna is to radiate energy. If we take a charged body and accelerate it up and down it radiates energy; if it were not charged it would not radiate energy. It is one thing to calculate from the conservation of energy that energy is lost, but another thing to answer the question, *against what force* are we doing the work? That is an interesting and very difficult question which has never been completely and satisfactorily answered for electrons, although it has been for antennas. What happens is this: in an antenna, the fields produced by the moving charges in one part of the antenna react on the moving charges in another part of the antenna. We can calculate these forces and find out how much work they do, and so find the right rule for the radiation resistance. When we say “We can calculate—” that is not quite right—we cannot, because we have not yet studied the laws of electricity at short distances; only at large distances do we know what the electric field is. We saw the formula (28.3), but at present it is too complicated for us to calculate the fields inside the wave zone. Of course, since conservation of energy is valid, we can calculate the result all right without knowing the fields at short distances. (As a matter of fact, by using this argument backwards it turns out that one can find the formula for the forces at short

32-1 Radiation resistance

32-2 The rate of radiation of energy

32-3 Radiation damping

32-4 Independent sources

32-5 Scattering of light

distances only by knowing the field at very large distances, by using the laws of conservation of energy, but we shall not go into that here.)

The problem in the case of a single electron is this: if there is only one charge, what can the force act on? It has been proposed, in the old classical theory, that the charge was a little ball, and that one part of the charge acted on the other part. Because of the delay in the action across the tiny electron, the force is not exactly in phase with the motion. That is, if we have the electron standing still, we know that "action equals reaction." So the various internal forces are equal, and there is no net force. But if the electron is accelerating, then because of the time delay across it, the force which is acting on the front from the back is not exactly the same as the force on the back from the front, because of the delay in the effect. This delay in the timing makes for a lack of balance, so, as a net effect, the thing holds itself back by its bootstraps! This model of the origin of the resistance to acceleration, the radiation resistance of a moving charge, has run into many difficulties, because our present view of the electron is that it is *not* a "little ball"; this problem has never been solved. Nevertheless we can calculate exactly, of course, what the net radiation resistance force must be, i.e., how much loss there must be when we accelerate a charge, in spite of not knowing directly the mechanism of how that force works.

32-2 The rate of radiation of energy

Now we shall calculate the total energy radiated by an accelerating charge. To keep the discussion general, we shall take the case of a charge accelerating any which way, but nonrelativistically. At a moment when the acceleration is, say, vertical, we know that the electric field that is generated is the charge multiplied by the projection of the retarded acceleration, divided by the distance. So we know the electric field at any point, and we therefore know the square of the electric field and thus the energy $\epsilon_0 c E^2$ leaving through a unit area per second.

The quantity $\epsilon_0 c$ appears quite often in expressions involving radiowave propagation. Its reciprocal is called the *impedance of a vacuum*, and it is an easy number to remember: it has the value $1/\epsilon_0 c = 377$ ohms. So the power in watts per square meter is equal to the average of the field squared, divided by 377.

Using our expression (29.1) for the electric field, we find that

$$S = \frac{q^2 a'^2 \sin^2 \theta}{16\pi^2 \epsilon_0 r^2 c^3} \quad (32.2)$$

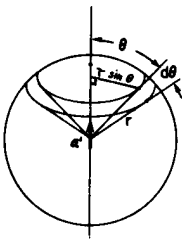


Fig. 32-1. The area of a spherical segment is $2\pi r \sin \theta \cdot r d\theta$.

is the power per square meter radiated in the direction θ . We notice that it goes inversely as the square of the distance, as we said before. Now suppose we wanted the total energy radiated in all directions: then we must integrate (32.2) over all directions. First we multiply by the area, to find the amount that flows within a little angle $d\theta$ (Fig. 32-1). We need the area of a spherical section. The way to think of it is this; if r is the radius, then the width of the annular segment is $r d\theta$, and the circumference is $2\pi r \sin \theta$, because $r \sin \theta$ is the radius of the circle. So the area of the little piece of the sphere is $2\pi r \sin \theta$ times $r d\theta$:

$$dA = 2\pi r^2 \sin \theta d\theta. \quad (32.3)$$

By multiplying the flux [(32.2), the power per square meter] by the area in square meters included in the small angle $d\theta$, we find the amount of energy that is liberated in this direction between θ and $\theta + d\theta$; then we integrate that over all the angles θ from 0 to 180° :

$$P = \int S dA = \frac{q^2 a'^2}{8\pi \epsilon_0 c^3} \int_0^\pi \sin^3 \theta d\theta. \quad (32.4)$$

By writing $\sin^3 \theta = (1 - \cos^2 \theta) \sin \theta$ it is not hard to show that $\int_0^\pi \sin^3 \theta d\theta = 4/3$. Using that fact, we finally get

$$P = \frac{q^2 a'^2}{6\pi \epsilon_0 c^3}. \quad (32.5)$$

This expression deserves some remarks. First of all, since the vector \mathbf{a}' had a certain direction, the a'^2 in (32.5) would be the square of the vector \mathbf{a}' , that is, $\mathbf{a}' \cdot \mathbf{a}'$, the length of the vector, squared. Secondly, the flux (32.2) was calculated using the retarded acceleration; that is, the acceleration at the time at which the energy now passing through the sphere was radiated. We might like to say that this energy was in fact liberated at this earlier time. This is not exactly true; it is only an approximate idea. The exact time when the energy is liberated cannot be defined precisely. All we can really calculate precisely is what happens in a complete motion, like an oscillation or something, where the acceleration finally ceases. Then what we find is that the total energy flux per cycle is the average of acceleration squared, for a complete cycle. This is what should really appear in (32.5). Or, if it is a motion with an acceleration that is initially and finally zero, then the total energy that has flown out is the time integral of (32.5).

To illustrate the consequences of formula (32.5) when we have an oscillating system, let us see what happens if the displacement x of the charge is oscillating so that the acceleration a is $-\omega^2 x_0 e^{i\omega t}$. The average of the acceleration squared over a cycle (remember that we have to be very careful when we square things that are written in complex notation—it really is the cosine, and the average of $\cos^2 \omega t$ is one-half) thus is

$$\langle a'^2 \rangle = \frac{1}{2} \omega^4 x_0^2.$$

Therefore

$$P = \frac{q^2 \omega^4 x_0^2}{12 \pi \epsilon_0 c^3}. \quad (32.6)$$

The formulas we are now discussing are relatively advanced and more or less modern; they date from the beginning of the twentieth century, and they are very famous. Because of their historical value, it is important for us to be able to read about them in older books. In fact, the older books also used a system of units different from our present mks system. However, all these complications can be straightened out in the final formulas dealing with electrons by the following rule: The quantity $q_e^2/4\pi\epsilon_0$, where q_e is the electronic charge (in coulombs), has, historically, been written as e^2 . It is very easy to calculate that e in the mks system is numerically equal to 1.5188×10^{-14} , because we know that, numerically, $q_e = 1.60206 \times 10^{-19}$ and $1/4\pi\epsilon_0 = 8.98748 \times 10^9$. Therefore we shall often use the convenient abbreviation

$$e^2 = \frac{q_e^2}{4\pi\epsilon_0}. \quad (32.7)$$

If we use the above numerical value of e in the older formulas and treat them as though they were written in mks units, we will get the right numerical results. For example, the older form of (32.5) is $P = \frac{2}{3} e^2 a^2 / c^3$. Again, the potential energy of a proton and an electron at distance r is $q_e^2/4\pi\epsilon_0 r$ or e^2/r , with $e = 1.5188 \times 10^{-14}$ mks.

32-3 Radiation damping

Now the fact that an oscillator loses a certain energy would mean that if we had a charge on the end of a spring (or an electron in an atom) which has a natural frequency ω_0 , and we start it oscillating and let it go, it will not oscillate forever, even if it is in empty space millions of miles from anything. There is no oil, no resistance, in an ordinary sense; no “viscosity.” But nevertheless it will not oscillate, as we might once have said, “forever,” because if it is charged it is radiating energy, and therefore the oscillation will slowly die out. How slowly? What is the Q of such an oscillator, caused by the electromagnetic effects, the so-called radiation resistance or radiation damping of the oscillator? The Q of any oscillating system is the total energy content of the oscillator at any time divided by the energy loss per radian:

$$Q = \frac{W}{dW/d\phi}.$$

Or (another way to write it), since $dW/d\phi = (dW/dt)/(d\phi/dt) = (dW/dt)/\omega$,

$$Q = \frac{\omega W}{dW/dt}. \quad (32.8)$$

If for a given Q this tells us how the energy of the oscillation dies out, $dW/dt = -(\omega/Q)W$, which has the solution $W = W_0 e^{-\omega t/Q}$ if W_0 is the initial energy (at $t = 0$).

To find the Q for a radiator, we go back to (32.8) and use (32.6) for dW/dt .

Now what do we use for the energy W of the oscillator? The kinetic energy of the oscillator is $\frac{1}{2}mv^2$, and the *mean* kinetic energy is $m\omega^2 x_0^2/4$. But we remember that for the total energy of an oscillator, on the average half is kinetic and half is potential energy, and so we double our result, and find for the total energy of the oscillator

$$W = \frac{1}{2}m\omega^2 x_0^2. \quad (32.9)$$

What do we use for the frequency in our formulas? We use the natural frequency ω_0 because, for all practical purposes, that is the frequency at which our atom is radiating, and for m we use the electron mass m_e . Then, making the necessary divisions and cancellations, the formula comes down to

$$\frac{1}{Q} = \frac{4\pi e^2}{3\lambda m_e c^2}. \quad (32.10)$$

(In order to see it better and in a more historical form we write it using our abbreviation $q_e^2/4\pi\epsilon_0 = e^2$, and the factor ω_0/c which was left over has been written as $2\pi/\lambda$.) Since Q is dimensionless, the combination $e^2/m_e c^2$ must be a property only of the electron charge and mass, an intrinsic property of the electron, and it must be a *length*. It has been given a name, the *classical electron radius*, because the early atomic models, which were invented to explain the radiation resistance on the basis of the force of one part of the electron acting on the other parts, all needed to have an electron whose dimensions were of this general order of magnitude. However, this quantity no longer has the significance that we believe that the electron really *has* such a radius. Numerically, the magnitude of the radius is

$$r_0 = \frac{e^2}{m_e c^2} = 2.82 \times 10^{-15} \text{ m}. \quad (32.11)$$

Now let us actually calculate the Q of an atom that is emitting light—let us say a sodium atom. For a sodium atom, the wavelength is roughly 6000 angstroms, in the yellow part of the visible spectrum, and this is a typical wavelength. Thus

$$Q = \frac{3\lambda}{4\pi r_0} \approx 5 \times 10^7, \quad (32.12)$$

so the Q of an atom is of the order 10^8 . This means that an atomic oscillator will oscillate for 10^8 radians or about 10^7 oscillations, before its energy falls by a factor $1/e$. The frequency of oscillation of light corresponding to 6000 angstroms, $\nu = c/\lambda$, is on the order of 10^{15} cycles/sec, and therefore the lifetime, the time it takes for the energy of a radiating atom to die out by a factor $1/e$, is on the order of 10^{-8} sec. In ordinary circumstances, freely emitting atoms usually take about this long to radiate. This is valid only for atoms which are in empty space, not being disturbed in any way. If the electron is in a solid and it has to hit other atoms or other electrons, then there are additional resistances and different damping.

The effective resistance term γ in the resistance law for the oscillator can be found from the relation $1/Q = \gamma/\omega_0$, and we remember that the size of γ determines how wide the resonance curve is (Fig. 23-2). Thus we have just computed the *widths of spectral lines* for freely radiating atoms! Since $\lambda = 2\pi c/\omega$, we find that

$$\begin{aligned} \Delta\lambda &= 2\pi c \Delta\omega/\omega^2 = 2\pi c \gamma/\omega_0^2 = 2\pi c/Q\omega_0 \\ &= \lambda/Q = 4\pi r_0/3 = 1.18 \times 10^{-14} \text{ m}. \end{aligned} \quad (32.13)$$

32-4 Independent sources

In preparation for our second topic, the scattering of light, we must now discuss a certain feature of the phenomenon of interference that we neglected to discuss previously. This is the question of when interference does *not* occur. If we have two sources S_1 and S_2 , with amplitudes A_1 and A_2 , and we make an observation in a certain direction in which the phases of arrival of the two signals are ϕ_1 and ϕ_2 (a combination of the actual time of oscillation and the delayed time, depending on the position of observation), then the energy that we receive can be found by compounding the two complex number vectors A_1 and A_2 , one at angle ϕ_1 and the other at angle ϕ_2 (as we did in Chapter 30) and we find that the resultant energy is proportional to

$$A_R^2 = A_1^2 + A_2^2 + 2A_1A_2 \cos(\phi_1 - \phi_2). \quad (32.14)$$

Now if the cross term $2A_1A_2 \cos(\phi_1 - \phi_2)$ were not there, then the total energy that would be received in a given direction would simply be the sum of the energies, $A_1^2 + A_2^2$, that would be liberated by each source separately, which is what we usually expect. That is, the combined intensity of light shining on something from two sources is the sum of the intensities of the two lights. On the other hand, if we have things set just right and we have a cross term, it is not such a sum, because there is also some interference. If there are circumstances in which this term is of no importance, then we would say the interference is apparently lost. Of course, in nature it is always there, but we may not be able to detect it.

Let us consider some examples. Suppose, first, that the two sources are 7,000,000,000 wavelengths apart, not an impossible arrangement. Then in a given direction it is true that there is a very definite value of these phase differences. But, on the other hand, if we move just a hair in one direction, a few wavelengths, which is no distance at all (our eye already has a hole in it that is so large that we are averaging the effects over a range very wide compared with one wavelength) then we change the relative phase, and the cosine changes very rapidly. If we take the *average* of the intensity over a little region, then the cosine, which goes plus, minus, plus, minus, as we move around, averages to zero.

So if we average over regions where the phase varies very rapidly with position, we get no interference.

Another example. Suppose that the two sources are two independent radio oscillators—not a single oscillator being fed by two wires, which guarantees that the phases are kept together, but two independent sources—and that they are not *precisely* tuned at the same frequency (it is very hard to make them at exactly the same frequency without actually wiring them together). In this case we have what we call two *independent* sources. Of course, since the frequencies are not exactly equal, although they started in phase, one of them begins to get a little ahead of the other, and pretty soon they are out of phase, and then it gets still further ahead, and pretty soon they are in phase again. So the phase difference between the two is gradually drifting with time, but if our observation is so crude that we cannot see that little time, if we average over a much longer time, then although the intensity swells and falls like what we call “beats” in sound, if these swellings and fallings are too rapid for our equipment to follow, then again this term averages out.

In other words, in any circumstance in which the phase shift averages out, we get no interference!

One finds many books which say that two distinct light sources never interfere. This is not a statement of physics, but is merely a statement of the degree of sensitivity of the technique of the experiments at the time the book was written. What happens in a light source is that first one atom radiates, then another atom radiates, and so forth, and we have just seen that atoms radiate a train of waves only for about 10^{-8} sec; after 10^{-8} sec, some atom has probably taken over, then another atom takes over, and so on. So the phases can really only stay the same for about 10^{-8} sec. Therefore, if we average for very much more than 10^{-8} sec, we do not see an interference from two different sources, because they cannot hold their phases steady for longer than 10^{-8} sec. With photocells, very high-speed detection

is possible, and one can show that there is an interference which varies with time, up and down, in about 10^{-8} sec. But most detection equipment, of course, does not look at such fine time intervals, and thus sees no interference. Certainly with the eye, which has a tenth-of-a-second averaging time, there is no chance whatever of seeing an interference between two different ordinary sources.

Recently it has become possible to make light sources which get around this effect by making all the atoms emit *together* in time. The device which does this is a very complicated thing, and has to be understood in a quantum-mechanical way. It is called a *laser*, and it is possible to produce from a laser a source in which the interference frequency, the time at which the phase is kept constant, is very much longer than 10^{-8} sec. It can be of the order of a hundredth, a tenth, or even one second, and so, with ordinary photocells, one can pick up the frequency between two different lasers. One can easily detect the pulsing of the beats between two laser sources. Soon, no doubt, someone will be able to demonstrate two sources shining on a wall, in which the beats are so slow that one can *see* the wall get bright and dark!

Another case in which the interference averages out is that in which, instead of having only *two* sources, we have *many*. In this case, we would write the expression for A_R^2 as the sum of a whole lot of amplitudes, complex numbers, squared, and we would get the square of each one, all added together, plus cross terms between every pair, and if the circumstances are such that the latter average out, then there will be no effects of interference. It may be that the various sources are located in such random positions that, although the phase difference between A_2 and A_3 is also definite, it is very different from that between A_1 and A_2 , etc. So we would get a whole lot of cosines, many plus, many minus, all averaging out.

So it is that in many circumstances we do not see the effects of interference, but see only a collective, total intensity equal to the sum of all the intensities.

32-5 Scattering of light

The above leads us to an effect which occurs in air as a consequence of the irregular positions of the atoms. When we were discussing the index of refraction, we saw that an incoming beam of light will make the atoms radiate again. The electric field of the incoming beam drives the electrons up and down, and they radiate because of their acceleration. This scattered radiation combines to give a beam in the same direction as the incoming beam, but of somewhat different phase, and this is the origin of the index of refraction.

But what can we say about the amount of re-radiated light in some other direction? Ordinarily, if the atoms are very beautifully located in a nice pattern, it is easy to show that we get nothing in other directions, because we are adding a lot of vectors with their phases always changing, and the result comes to zero. But if the objects are *randomly* located, then the total intensity in any direction is the *sum* of the intensities that are scattered by each atom, as we have just discussed. Furthermore, the atoms in a gas are in actual motion, so that although the relative phase of two atoms is a definite amount now, later the phase would be quite different, and therefore *each* cosine term will average out. Therefore, to find out how much light is scattered in a given direction by a gas, we merely study the effects of *one atom* and multiply the intensity it radiates by the number of atoms.

Earlier, we remarked that the phenomenon of scattering of light of this nature is the origin of the blue of the sky. The sunlight goes through the air, and when we look to one side of the sun—say at 90° to the beam—we see blue light; what we now have to calculate is *how much* light we see and *why it is blue*.

If the incident beam has the electric field $\mathbf{E} = \mathbf{E}_0 e^{i\omega t}$ at the point where the atom is located, we know that an electron in the atom will vibrate up and down in response to this \mathbf{E} (Fig. 32-2). From Eq. (23.8), the amplitude will be

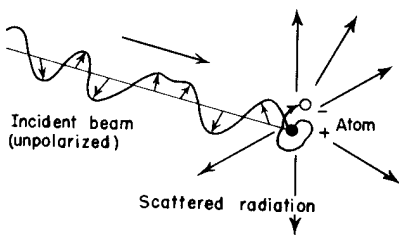


Fig. 32-2. A beam of radiation falls on an atom and causes the charges (electrons) in the atom to move. The moving electrons in turn radiate in various directions.

$$\hat{x} = \frac{q_e \mathbf{E}_0}{m(\omega_0^2 - \omega^2 + i\omega\gamma)}. \quad (32.15)$$

We could include the damping and the possibility that the atom acts like several oscillators of different frequency and sum over the various frequencies, but for simplicity let us just take one oscillator and neglect the damping. Then the response to the external electric field, which we have already used in the calculation of the index of refraction, is simply

$$\hat{\mathbf{x}} = \frac{q_e \hat{\mathbf{E}}_0}{m(\omega_0^2 - \omega^2)}. \quad (32.16)$$

We could now easily calculate the intensity of light that is emitted in various directions, using formula (32.2) and the acceleration corresponding to the above $\hat{\mathbf{x}}$.

Rather than do this, however, we shall simply calculate the *total amount* of light scattered in *all* directions, just to save time. The total amount of light energy per second, scattered in all directions by the single atom, is of course given by Eq. (32.7). So, putting together the various pieces and regrouping them, we get

$$\begin{aligned} P &= [(q_e^2 \omega^4 / 12 \pi \epsilon_0 c^3) q_e^2 E_0^2 / m^2 (\omega^2 - \omega_0^2)^2] \\ &= (\frac{1}{2} \epsilon_0 c E_0^2) (8 \pi / 3) (q_e^4 / 16 \pi^2 \epsilon_0^2 m^2 c^4) [\omega^4 / (\omega^2 - \omega_0^2)^2] \\ &= (\frac{1}{2} \epsilon_0 c E_0^2) (8 \pi r_0^2 / 3) [\omega^4 / (\omega^2 - \omega_0^2)^2] \end{aligned} \quad (32.17)$$

for the total scattered power, radiated in all directions.

We have written the result in the above form because it is then easy to remember: First, the total energy that is scattered is proportional to the square of the incident field. What does that mean? Obviously, the square of the incident field is proportional to the energy which is coming in per second. In fact, the energy incident per square meter per second is $\epsilon_0 c$ times the average $\langle E^2 \rangle$ of the square of the electric field, and if E_0 is the maximum value of E , then $\langle E^2 \rangle = \frac{1}{2} E_0^2$. In other words, the total energy scattered is proportional to the energy per square meter that comes in; the brighter the sunlight that is shining in the sky, the brighter the sky is going to look.

Next, what *fraction* of the incoming light is scattered? Let us imagine a “target” with a certain area, let us say σ , in the beam (not a real, material target, because this would diffract light, and so on; we mean an imaginary area drawn in space). The total amount of energy that would pass through this surface σ in a given circumstance is proportional both to the incoming intensity and to σ , and would be

$$P = (\frac{1}{2} \epsilon_0 c E_0^2) \sigma. \quad (32.18)$$

Now we invent an idea: we say that the atom scatters a total amount of intensity which is the amount which would fall on a certain geometrical area, and we give the answer by giving that area. That answer, then, is independent of the incident intensity; it gives the ratio of the energy scattered to the energy incident per square meter. In other words, the ratio

$$\frac{\text{total energy scattered per second}}{\text{energy incident per square meter per second}} \text{ is an } \textit{area}.$$

The significance of this area is that, if all the energy that impinged on that area were to be spewed in all directions, then that is the amount of energy that would be scattered by the atom.

This area is called a *cross section for scattering*; the idea of cross section is used constantly, whenever some phenomenon occurs in proportion to the intensity of a beam. In such cases one always describes the amount of the phenomenon by saying what the effective area would have to be to pick up that much of the beam. It does not mean in any way that this oscillator actually *has* such an area. If there were nothing present but a free electron shaking up and down there would be no area directly associated with it, physically. It is merely a way of expressing the answer to a certain kind of problem; it tells us what area the incident beam would

have to hit in order to account for that much energy coming off. Thus, for our case,

$$\sigma_s = \frac{8\pi r_0^2}{3} \cdot \frac{\omega^4}{(\omega^2 - \omega_0^2)^2} \quad (32.19)$$

(the subscript s is for “scattering”).

Let us look at some examples. First, if we go to a very low natural frequency ω_0 , or to completely unbound electrons, for which $\omega_0 = 0$, then the frequency ω cancels out and the cross section is a constant. This low-frequency limit, or the free electron cross section, is known as the *Thompson scattering cross section*. It is an area whose dimensions are approximately 10^{-15} meter, more or less, on a side, i.e., 10^{-30} square meter, which is rather small!

On the other hand, if we take the case of light in the air, we remember that for air the natural frequencies of the oscillators are higher than the frequency of the light that we use. This means that, to a first approximation, we can disregard ω^2 in the denominator, and we find that the scattering is proportional to the *fourth power* of the frequency. That is to say, light which is of higher frequency by, say, a factor of two, is *sixteen times* more intensely scattered, which is a quite sizable difference. This means that blue light, which has about twice the frequency of the reddish end of the spectrum, is scattered to a far greater extent than red light. Thus when we look at the sky it looks that glorious blue that we see all the time!

There are several points to be made about the above results. One interesting question is, why do we ever see the *clouds*? Where do the clouds come from? Everybody knows it is the condensation of water vapor. But, of course, the water vapor is already in the atmosphere *before* it condenses, so why don't we see it then? After it condenses it is perfectly obvious. It wasn't there, now it *is* there. So the mystery of where the clouds come from is not really such a childish mystery as “Where does the water come from, Daddy?,” but has to be explained.

We have just explained that every atom scatters light, and of course the water vapor will scatter light, too. The mystery is why, when the water is condensed into clouds, does it scatter such a *tremendously greater amount* of light?

Consider what would happen if, instead of a single atom, we had an agglomerate of atoms, say two, very close together compared with the wavelength of the light. Remember, atoms are only an angstrom or so across, while the wavelength of light is some 5000 angstroms, so when they form a clump, a few atoms together, they can *be* very close together compared with the wavelength of light. Then when the electric field acts, *both of the atoms will move together*. The electric field that is scattered will then be the sum of the two electric fields in phase, i.e., double the amplitude that there was with a single atom, and the *energy* which is scattered is therefore *four times* what it is with a single atom, not twice! So lumps of atoms radiate or scatter more energy than they do as single atoms. Our argument that the phases are independent is based on the assumption that there is a real and large difference in phase between any two atoms, which is true only if they are several wavelengths apart and randomly spaced, or moving. But if they are right next to each other, they necessarily scatter in phase, and they have a coherent interference which produces an increase in the scattering.

If we have N atoms in a lump, which is a tiny droplet of water, then each one will be driven by the electric field in about the same way as before (the effect of one atom on the other is not important; it is just to get the idea anyway) and the amplitude of scattering from each one is the same, so the total field which is scattered is N -fold increased. The *intensity* of the light which is scattered is then the square, or N^2 -fold, increased. We would have expected, if the atoms were spread out in space, only N times as much as 1, whereas we get N^2 times as much as 1! That is to say, the scattering of water in lumps of N molecules each is N times more intense than the scattering of the single atoms. So as the water agglomerates the scattering increases. Does it increase *ad infinitum*? No! When does this analysis begin to fail? How many atoms can we put together before we cannot drive this argument any further? *Answer*: If the water drop gets so big that from one end to the other is a wavelength or so, then the atoms are no longer all in

phase because they are too far apart. So as we keep increasing the size of the drop-lets we get more and more scattering, until such a time that a drop gets about the size of a wavelength, and then the scattering does not increase anywhere nearly as rapidly as the drop gets bigger. Furthermore, the blue disappears, because for long wavelengths the drops can be bigger, before this limit is reached, than they can be for short wavelengths. Although the short waves scatter more per atom than the long waves, there is a bigger enhancement for the red end of the spectrum than for the blue end when all the drops are bigger than the wavelength, so the color is shifted from the blue toward the red.

Now we can make an experiment that demonstrates this. We can make particles that are very small at first, and then gradually grow in size. We use a solution of sodium thiosulfate (hypo) with sulphuric acid, which precipitates very fine grains of sulphur. As the sulphur precipitates, the grains first start very small, and the scattering is a little bluish. As it precipitates more it gets more intense, and then it will get whitish as the particles get bigger. In addition, the light which goes straight through will have the blue taken out. That is why the sunset is red, of course, because the light that comes through a lot of air to the eye has had a lot of blue light scattered out, so it is yellow-red.

Finally, there is one other important feature which really belongs in the next chapter, on polarization, but it is so interesting that we point it out now. This is that the electric field of the scattered light tends to vibrate in a particular direction. The electric field in the incoming light is oscillating in some way, and the driven oscillator goes in this same direction, and if we are situated about at right angles to the beam, we will see *polarized* light, that is to say, light in which the electric field is going only one way. In general, the atoms can vibrate in any direction at right angles to the beam, but if they are driven directly toward or away from us, we do not see it. So if the incoming light has an electric field which changes and oscillates in any direction, which we call unpolarized light, then the light which is coming out at 90° to the beam vibrates in only one direction! (See Fig. 32-3.)

There is a substance called polaroid which has the property that when light goes through it, only the piece of the electric field which is along one particular axis can get through. We can use this to test for polarization, and indeed we find the light scattered by the hypo solution to be strongly polarized.

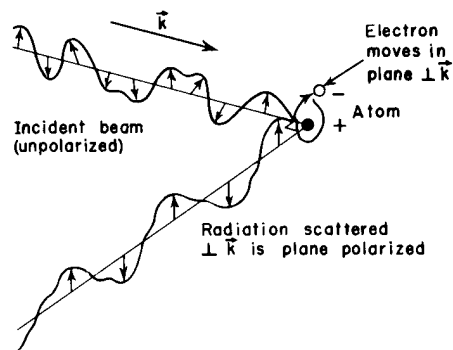


Fig. 32-3. Illustration of the origin of the polarization of radiation scattered at right angles to the incident beam.

Polarization

33-1 The electric vector of light

In this chapter we shall consider those phenomena which depend on the fact that the electric field that describes the light is a vector. In previous chapters we have not been concerned with the direction of oscillation of the electric field, except to note that the electric vector lies in a plane perpendicular to the direction of propagation. The particular direction in this plane has not concerned us. We now consider those phenomena whose central feature is the particular direction of oscillation of the electric field.

In ideally monochromatic light, the electric field must oscillate at a definite frequency, but since the x -component and the y -component can oscillate independently at a definite frequency, we must first consider the resultant effect produced by superposing two independent oscillations at right angles to each other. What kind of electric field is made up of an x -component and a y -component which oscillate at the same frequency? If one adds to an x -vibration a certain amount of y -vibration at the same phase, the result is a vibration in a new direction in the xy -plane. Figure 33-1 illustrates the superposition of different amplitudes for the x -vibration and the y -vibration. But the resultants shown in Fig. 33-1 are not the only possibilities; in all of these cases we have assumed that the x -vibration and the y -vibration are *in phase*, but it does not have to be that way. It could be that the x -vibration and the y -vibration are out of phase.

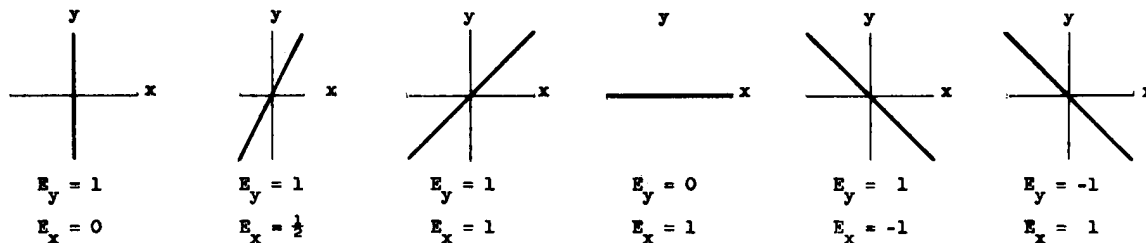


Fig. 33-1. Superposition of x -vibrations and y -vibrations in phase.

When the x -vibration and the y -vibration are not in phase, the electric field vector moves around in an ellipse, and we can illustrate this in a familiar way. If we hang a ball from a support by a long string, so that it can swing freely in a horizontal plane, it will execute sinusoidal oscillations. If we imagine horizontal x - and y -coordinates with their origin at the rest position of the ball, the ball can swing in either the x - or y -direction with the same pendulum frequency. By selecting the proper initial displacement and initial velocity, we can set the ball in oscillation along either the x -axis or the y -axis, or along any straight line in the xy -plane. These motions of the ball are analogous to the oscillations of the electric field vector illustrated in Fig. 33-1. In each instance, since the x -vibrations and the y -vibrations reach their maxima and minima at the same time, the x - and y -oscillations are in phase. But we know that the most general motion of the ball is motion in an ellipse, which corresponds to oscillations in which the x - and y -directions are *not* in the same phase. The superposition of x - and y -vibrations which are not in phase is illustrated in Fig. 33-2 for a variety of angles between the phase of the x -vibration and that of the y -vibration. The general result is that the electric vector moves around an ellipse. The motion in a straight line is a particular

33-1 The electric vector of light

33-2 Polarization of scattered light

33-3 Birefringence

33-4 Polarizers

33-5 Optical activity

33-6 The intensity of reflected light

33-7 Anomalous refraction

case corresponding to a phase difference of zero (or an integral multiple of π); motion in a circle corresponds to equal amplitudes with a phase difference of 90° (or any odd integral multiple of $\pi/2$).

In Fig. 33-2 we have labeled the electric field vectors in the x - and y -directions with complex numbers, which are a convenient representation in which to express the phase difference. Do not confuse the real and imaginary components of the complex electric vector in this notation with the x - and y -coordinates of the field. The x - and y -coordinates plotted in Fig. 33-1 and Fig. 33-2 are actual electric fields that we can measure. The real and imaginary components of a complex electric field vector are only a mathematical convenience and have no physical significance.

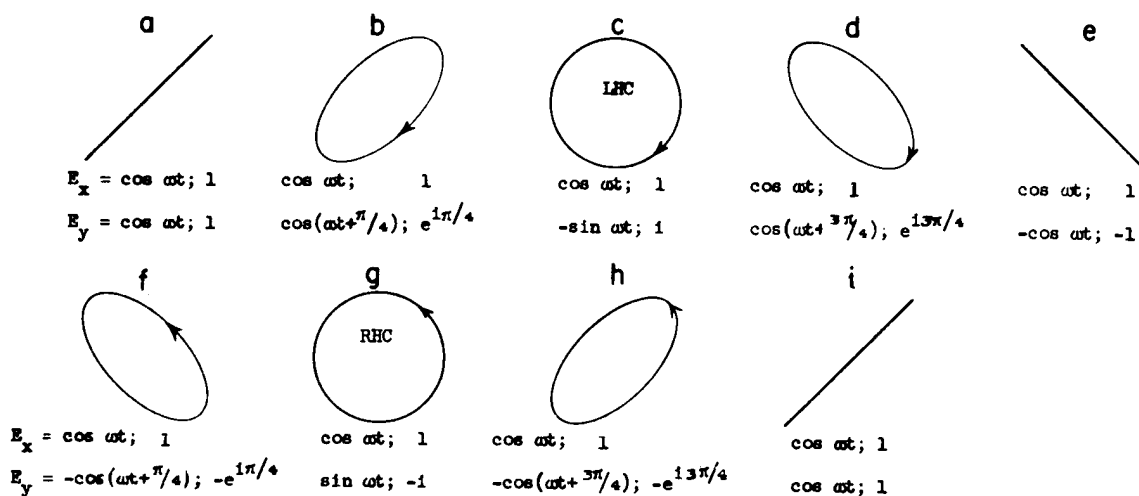


Fig. 33-2. Superposition of x -vibrations and y -vibrations with equal amplitudes but various relative phases. The components E_x and E_y are expressed in both real and complex notations.

Now for some terminology. Light is *linearly polarized* (sometimes called plane polarized) when the electric field oscillates on a straight line; Fig. 33-1 illustrates linear polarization. When the end of the electric field vector travels in an ellipse, the light is *elliptically polarized*. When the end of the electric field vector travels around a circle, we have *circular polarization*. If the end of the electric vector, when we look at it as the light comes straight toward us, goes around in a counterclockwise direction, we call it right-hand circular polarization. Figure 33-2(g) illustrates right-hand circular polarization, and Fig. 33-2(c) shows left-hand circular polarization. In both cases the light is coming out of the paper. Our convention for labeling left-hand and right-hand circular polarization is consistent with that which is used today for all the other particles in physics which exhibit polarization (e.g., electrons). However, in some books on optics the opposite conventions are used, so one must be careful.

We have considered linearly, circularly, and elliptically polarized light, which covers everything except for the case of *unpolarized* light. Now how can the light be unpolarized when we know that it must vibrate in one or another of these ellipses? If the light is not absolutely monochromatic, or if the x - and y -phases are not kept perfectly together, so that the electric vector first vibrates in one direction, then in another, the polarization is constantly changing. Remember that one atom emits during 10^{-8} sec, and if one atom emits a certain polarization, and then another atom emits light with a different polarization, the polarizations will change every 10^{-8} sec. If the polarization changes more rapidly than we can detect it, then we call the light unpolarized, because all the effects of the polarization average out. None of the interference effects of polarization would show up with unpolarized light. But as we see from the definition, light is unpolarized only if we are unable to find out whether the light is polarized or not.

33-2 Polarization of scattered light

The first example of the polarization effect that we have already discussed is the scattering of light. Consider a beam of light, for example from the sun, shining on the air. The electric field will produce oscillations of charges in the air, and motion of these charges will radiate light with its maximum intensity in a plane normal to the direction of vibration of the charges. The beam from the sun is unpolarized, so the direction of polarization changes constantly, and the direction of vibration of the charges in the air changes constantly. If we consider light scattered at 90° , the vibration of the charged particles radiates to the observer only when the vibration is perpendicular to the observer's line of sight, and then light will be polarized along the direction of vibration. So scattering is an example of one means of producing polarization.

33-3 Birefringence

Another interesting effect of polarization is the fact that there are substances for which the index of refraction is different for light linearly polarized in one direction and linearly polarized in another. Suppose that we had some material which consisted of long, nonspherical molecules, longer than they are wide, and suppose that these molecules were arranged in the substance with their long axes parallel. Then what happens when the oscillating electric field passes through this substance? Suppose that because of the structure of the molecule, the electrons in the substance respond more easily to oscillations in the direction parallel to the axes of the molecules than they would respond if the electric field tries to push them at right angles to the molecular axis. In this way we expect a different response for polarization in one direction than for polarization at right angles to that direction. Let us call the direction of the axes of the molecules the *optic axis*. When the polarization is in the direction of the optic axis the index of refraction is different than it would be if the direction of polarization were at right angles to it. Such a substance is called *birefringent*. It has two refrangibilities, i.e., two indexes of refraction, depending on the direction of the polarization inside the substance. What kind of a substance can be birefringent? In a birefringent substance there must be a certain amount of lining up, for one reason or another, of unsymmetrical molecules. Certainly a cubic crystal, which has the symmetry of a cube, cannot be birefringent. But long needlelike crystals undoubtedly contain molecules that are asymmetric, and one observes this effect very easily.

Let us see what effects we would expect if we were to shine polarized light through a plate of a birefringent substance. If the polarization is parallel to the optic axis, the light will go through with one velocity; if the polarization is perpendicular to the axis, the light is transmitted with a different velocity. An interesting situation arises when, say, light is linearly polarized at 45° to the optic axis. Now the 45° polarization, we have already noticed, can be represented as a superposition of the x - and the y -polarizations of equal amplitude and in phase, as shown in Fig. 33-2(a). Since the x - and y -polarizations travel with different velocities, their phases change at a different rate as the light passes through the substance. So, although at the start the x - and y -vibrations are in phase, inside the material the phase difference between x - and y -vibrations is proportional to the depth in the substance. As the light proceeds through the material the polarization changes as shown in the series of diagrams in Fig. 33-2. If the thickness of the plate is just right to introduce a 90° phase shift between the x - and y -polarizations, as in Fig. 33-2(c), the light will come out circularly polarized. Such a thickness is called a quarter-wave plate, because it introduces a quarter-cycle phase difference between the x - and the y -polarizations. If linearly polarized light is sent through two quarter-wave plates, it will come out plane-polarized again, but at right angles to the original direction, as we can see from Fig. 33-2(e).

One can easily illustrate this phenomenon with a piece of cellophane. Cellophane is made of long, fibrous molecules, and is not isotropic, since the fibers lie preferentially in a certain direction. To demonstrate birefringence we need a

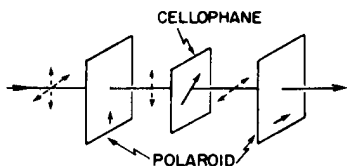


Fig. 33-3. An experimental demonstration of the birefringence of cellophane. The electric vectors in the light are indicated by the dotted lines. The pass axes of the polaroid sheets and optic axes of the cellophane are indicated by arrows. The incident beam is unpolarized.

beam of linearly polarized light, and we can obtain this conveniently by passing unpolarized light through a sheet of polaroid. Polaroid, which we will discuss later in more detail, has the useful property that it transmits light that is linearly polarized parallel to the axis of the polaroid with very little absorption, but light polarized in a direction perpendicular to the axis of the polaroid is strongly absorbed. When we pass unpolarized light through a sheet of polaroid, only that part of the unpolarized beam which is vibrating parallel to the axis of the polaroid gets through, so that the transmitted beam is linearly polarized. This same property of polaroid is also useful in detecting the direction of polarization of a linearly polarized beam, or in determining whether a beam is linearly polarized or not. One simply passes the beam of light through the polaroid sheet and rotates the polaroid in the plane normal to the beam. If the beam is linearly polarized, it will not be transmitted through the sheet when the axis of the polaroid is normal to the direction of polarization. The transmitted beam is only slightly attenuated when the axis of the polaroid sheet is rotated through 90° . If the transmitted intensity is independent of the orientation of the polaroid, the beam is not linearly polarized.

To demonstrate the birefringence of cellophane, we use two sheets of polaroid, as shown in Fig. 33-3. The first gives us a linearly polarized beam which we pass through the cellophane and then through the second polaroid sheet, which serves to detect any effect the cellophane may have had on the polarized light passing through it. If we first set the axes of the two polaroid sheets perpendicular to each other and remove the cellophane, no light will be transmitted through the second polaroid. If we now introduce the cellophane between the two polaroid sheets, and rotate the sheet about the beam axis, we observe that in general the cellophane makes it possible for some light to pass through the second polaroid. However, there are two orientations of the cellophane sheet, at right angles to each other, which permit no light to pass through the second polaroid. These orientations in which linearly polarized light is transmitted through the cellophane with no effect on the direction of polarization must be the directions parallel and perpendicular to the optic axis of the cellophane sheet.

We suppose that the light passes through the cellophane with two different velocities in these two different orientations, but it is transmitted without changing the direction of polarization. When the cellophane is turned halfway between these two orientations, as shown in Fig. 33-3, we see that the light transmitted through the second polaroid is bright.

It just happens that ordinary cellophane used in commercial packaging is very close to a half-wave thickness for most of the colors in white light. Such a sheet will turn the axis of linearly polarized light through 90° if the incident linearly polarized beam makes an angle of 45° with the optic axis, so that the beam emerging from the cellophane is then vibrating in the right direction to pass through the second polaroid sheet.

If we use white light in our demonstration, the cellophane sheet will be of the proper half-wave thickness only for a particular component of the white light, and the transmitted beam will have the color of this component. The color transmitted depends on the thickness of the cellophane sheet, and we can vary the effective thickness of the cellophane by tilting it so that the light passes through the cellophane at an angle, consequently through a longer path in the cellophane. As the sheet is tilted the transmitted color changes. With cellophane of different thicknesses one can construct filters that will transmit different colors. These filters have the interesting property that they transmit one color when the two polaroid sheets have their axes perpendicular, and the complementary color when the axes of the two polaroid sheets are parallel.

Another interesting application of aligned molecules is quite practical. Certain plastics are composed of very long and complicated molecules all twisted together. When the plastic is solidified very carefully, the molecules are all twisted in a mass, so that there are as many aligned in one direction as another, and so the plastic is not particularly birefringent. Usually there are strains and stresses introduced when the material is solidified, so the material is not perfectly homo-

geneous. However, if we apply tension to a piece of this plastic material, it is as if we were pulling a whole tangle of strings, and there will be more strings preferentially aligned parallel to the tension than in any other direction. So when a stress is applied to certain plastics, they become birefringent, and one can see the effects of the birefringence by passing polarized light through the plastic. If we examine the transmitted light through a polaroid sheet, patterns of light and dark fringes will be observed (in color, if white light is used). The patterns move as stress is applied to the sample, and by counting the fringes and seeing where most of them are, one can determine what the stress is. Engineers use this phenomenon as a means of finding the stresses in odd-shaped pieces that are difficult to calculate.

Another interesting example of a way of obtaining birefringence is by means of a liquid substance. Consider a liquid composed of long asymmetric molecules which carry a plus or minus average charge near the ends of the molecule, so that the molecule is an electric dipole. In the collisions in the liquid the molecules will ordinarily be randomly oriented, with as many molecules pointed in one direction as in another. If we apply an electric field the molecules will tend to line up, and the moment they line up the liquid becomes birefringent. With two polaroid sheets and a transparent cell containing such a polar liquid, we can devise an arrangement with the property that light is transmitted only when the electric field is applied. So we have an electrical switch for light, which is called a *Kerr cell*. This effect, that an electric field can produce birefringence in certain liquids, is called the Kerr effect.

33-4 Polarizers

So far we have considered substances in which the refractive index is different for light polarized in different directions. Of very practical value are those crystals and other substances in which not only the index, but also the coefficient of absorption, is different for light polarized in different directions. By the same arguments which supported the idea of birefringence, it is understandable that absorption can vary with the direction in which the charges are forced to vibrate in an anisotropic substance. Tourmaline is an old, famous example and polaroid is another. Polaroid consists of a thin layer of small crystals of herapathite (a salt of iodine and quinine), all aligned with their axes parallel. These crystals absorb light when the oscillations are in one direction, and they do not absorb appreciably when the oscillations are in the other direction.

Suppose that we send light into a polaroid sheet polarized linearly at an angle θ to the passing direction. What intensity will come through? This incident light can be resolved into a component perpendicular to the pass direction which is proportional to $\sin \theta$, and a component along the pass direction which is proportional to $\cos \theta$. The amplitude which comes out of the polaroid is only the cosine θ part; the $\sin \theta$ component is absorbed. The amplitude which passes through the polaroid is smaller than the amplitude which entered, by a factor $\cos \theta$. The energy which passes through the polaroid, i.e., the intensity of the light, is proportional to the square of $\cos \theta$. $\cos^2 \theta$, then, is the intensity transmitted when the light enters polarized at an angle θ to the pass direction. The absorbed intensity, of course, is $\sin^2 \theta$.

An interesting paradox is presented by the following situation. We know that it is not possible to send a beam of light through two polaroid sheets with their axes crossed at right angles. But if we place a third polaroid sheet *between* the first two, with its pass axis at 45° to the crossed axes, some light is transmitted. We know that polaroid absorbs light, it does not create anything. Nevertheless, the addition of a third polaroid at 45° allows more light to get through. The analysis of this phenomenon is left as an exercise for the student.

One of the most interesting examples of polarization is not in complicated crystals or difficult substances, but in one of the simplest and most familiar of situations—the reflection of light from a surface. Believe it or not, when light is reflected from a glass surface it may be polarized, and the physical explanation of this is very simple. It was discovered empirically by Brewster that light reflected

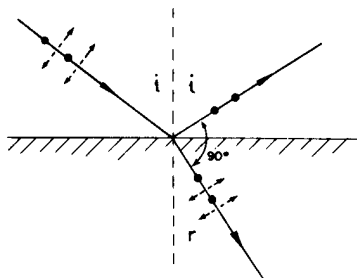


Fig. 33-4. Reflection of linearly polarized light at Brewster's angle. The polarization direction is indicated by dashed arrows; round dots indicate polarization normal to the paper.

from a surface is completely polarized if the reflected beam and the beam refracted into the material form a right angle. The situation is illustrated in Fig. 33-4. If the incident beam is polarized in the plane of incidence, there will be no reflection at all. Only if the incident beam is polarized normal to the plane of incidence will it be reflected. The reason is very easy to understand. In the reflecting material the light is polarized transversely, and we know that it is the motion of the charges in the material which generates the emergent beam, which we call the reflected beam. The source of this so-called reflected light is not simply that the incident beam is reflected; our deeper understanding of this phenomenon tells us that the incident beam drives an oscillation of the charges in the material, which in turn generates the reflected beam. From Fig. 33-4 it is clear that only oscillations normal to the paper can radiate in the direction of reflection, and consequently the reflected beam will be polarized normal to the plane of incidence. If the incident beam is polarized in the plane of incidence, there will be no reflected light.

This phenomenon is readily demonstrated by reflecting a linearly polarized beam from a flat piece of glass. If the glass is turned to present different angles of incidence to the polarized beam, sharp attenuation of the reflected intensity is observed when the angle of incidence passes through Brewster's angle. This attenuation is observed only if the plane of polarization lies in the plane of incidence. If the plane of polarization is normal to the plane of incidence, the usual reflected intensity is observed at all angles.

33-5 Optical activity

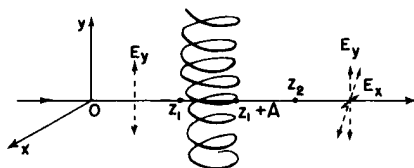


Fig. 33-5. A molecule with a shape that is not symmetric when reflected in a mirror. A beam of light, linearly polarized in the y -direction, falls on the molecule.

Another most remarkable effect of polarization is observed in materials composed of molecules which do not have reflection symmetry: molecules shaped something like a corkscrew, or like a gloved hand, or any shape which, if viewed through a mirror, would be reversed in the same way that a left-hand glove reflects as a right-hand glove. Suppose all of the molecules in the substance are the same, i.e., none is a mirror image of any other. Such a substance may show an interesting effect called optical activity, whereby as linearly polarized light passes through the substance, the direction of polarization rotates about the beam axis.

To understand the phenomenon of optical activity requires some calculation, but we can see qualitatively how the effect might come about, without actually carrying out the calculations. Consider an asymmetric molecule in the shape of a spiral, as shown in Fig. 33-5. Molecules need not actually be shaped like a corkscrew in order to exhibit optical activity, but this is a simple shape which we shall take as a typical example of those that do not have reflection symmetry. When a light beam linearly polarized along the y -direction falls on this molecule, the electric field will drive charges up and down the helix, thereby generating a current in the y -direction and radiating an electric field E_y polarized in the y -direction. However, if the electrons are constrained to move along the spiral, they must also move in the x -direction as they are driven up and down. When a current is flowing up the spiral, it is also flowing into the paper at $z = z_1$ and out of the paper at $z = z_1 + A$, if A is the diameter of our molecular spiral. One might suppose that the current in the x -direction would produce no net radiation, since the currents are in opposite directions on opposite sides of the spiral. However, if we consider the x -components of the electric field arriving at $z = z_2$, we see that the field radiated by the current at $z = z_1 + A$ and the field radiated from $z = z_1$ arrive at z_2 separated in time by the amount A/c , and thus separated in phase by $\pi + \omega A/c$. Since the phase difference is not exactly π , the two fields do not cancel exactly, and we are left with a small x -component in the electric field generated by the motion of the electrons in the molecule, whereas the driving electric field had only a y -component. This small x -component, added to the large y -component, produces a resultant field that is tilted slightly with respect to the y -axis, the original direction of polarization. As the light moves through the material, the direction of polarization rotates about the beam axis. By drawing a few examples and considering the currents that will be set in motion by an incident

electric field, one can convince himself that the existence of optical activity and the sign of the rotation are independent of the orientation of the molecules.

Corn syrup is a common substance which possesses optical activity. The phenomenon is easily demonstrated with a polaroid sheet to produce a linearly polarized beam, a transmission cell containing corn syrup, and a second polaroid sheet to detect the rotation of the direction of polarization as the light passes through the corn syrup.

33-6 The intensity of reflected light

Let us now consider quantitatively the reflection coefficient as a function of angle. Figure 33-6(a) shows a beam of light striking a glass surface, where it is partly reflected and partly refracted into the glass. Let us suppose that the incident beam, of unit amplitude, is linearly polarized normal to the plane of the paper. We will call the amplitude of the reflected wave b , and the amplitude of the refracted wave a . The refracted and reflected waves will, of course, be linearly polarized, and the electric field vectors of the incident, reflected, and refracted waves are all parallel to each other. Figure 33-6(b) shows the same situation, but now we suppose that the incident wave, of unit amplitude, is polarized in the plane of the paper. Now let us call the amplitude of the reflected and refracted wave B and A , respectively.

We wish to calculate how strong the reflection is in the two situations illustrated in Fig. 33-6(a) and 33-6(b). We already know that when the angle between the reflected beam and refracted beam is a right angle, there will be no reflected wave in Fig. 33-6(b), but let us see if we cannot get a quantitative answer—an exact formula for B and b as a function of the angle of incidence, i .

The principle that we must understand is as follows. The currents that are generated in the glass produce two waves. First, they produce the reflected wave. Moreover, we know that if there were no currents generated in the glass, the incident wave would continue straight into the glass. Remember that all the sources in the world make the net field. The source of the incident light beam produces a field of unit amplitude which would move into the glass along the dotted line in the figure. This field is not observed, and therefore the currents generated in the glass must produce a field of amplitude -1 which moves along the dotted line. Using this fact, we will calculate the amplitude of the refracted waves, a and A .

In Fig. 33-6(a) we see that the field of amplitude b is radiated by the motion of charges inside the glass which are responding to a field a inside the glass, and that therefore b is proportional to a . We might suppose that since our two figures are exactly the same, except for the direction of polarization, the ratio B/A would be the same as the ratio b/a . This is not quite true, however, because in Fig. 33-6(b) the polarization directions are not all parallel to each other, as they are in Fig. 33-6(a). It is only the component of A which is perpendicular to B , $A \cos(i + r)$, which is effective in producing B . The correct expression for the proportionality is then

$$\frac{b}{a} = \frac{B}{A \cos(i + r)}. \quad (33.1)$$

Now we use a trick. We know that in both (a) and (b) of Fig. 33-6 the electric field in the glass must produce oscillations of the charges which generate a field of amplitude -1 , polarized parallel to the incident beam, and moving in the direction of the dotted line. But we see from part (b) of the figure that only the component of A that is normal to the dashed line has the right polarization to produce this field, whereas in Fig. 33-6(a) the full amplitude a is effective, since the polarization of wave a is parallel to the polarization of the wave of amplitude -1 . Therefore we can write

$$\frac{A \cos(i - r)}{a} = \frac{-1}{-1}, \quad (33.2)$$

since the two amplitudes on the left side of Eq. (33.2) each produce the wave of amplitude -1 .

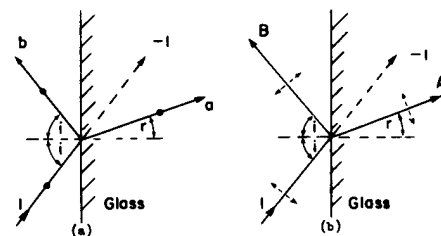


Fig. 33-6. An incident wave of unit amplitude is reflected and refracted at a glass surface. In (a) the incident wave is linearly polarized normal to the plane of the paper. In (b) the incident wave is linearly polarized in the direction shown by the dotted electric vector.

Dividing Eq. (33.1) by Eq. (33.2), we obtain

$$\frac{B}{b} = \frac{\cos(i+r)}{\cos(i-r)}, \quad (33.3)$$

a result which we can check against what we already know. If we set $(i+r) = 90^\circ$, Eq. (33.3) gives $B = 0$, as Brewster says it should be, so our results so far are at least not obviously wrong.

We have assumed unit amplitudes for the incident waves, so that $|B|^2/1^2$ is the reflection coefficient for waves polarized in the plane of incidence, and $|b|^2/1^2$ is the reflection coefficient for waves polarized normal to the plane of incidence. The ratio of these two reflection coefficients is determined by Eq. (33.3).

Now we perform a miracle, and compute not just the ratio, but each coefficient $|B|^2$ and $|b|^2$ individually! We know from the conservation of energy that the energy in the refracted wave must be equal to the incident energy minus the energy in the reflected wave, $1 - |B|^2$ in one case, $1 - |b|^2$ in the other. Furthermore, the energy which passes into the glass in Fig. 33-6(b) is to the energy which passes into the glass in Fig. 33-6(a) as the ratio of the squares of the refracted amplitudes, $|A|^2/|a|^2$. One might ask whether we really know how to compute the energy inside the glass, because, after all, there are energies of motion of the atoms in addition to the energy in the electric field. But it is obvious that all of the various contributions to the total energy will be proportional to the square of the amplitude of the electric field. Therefore we can write

$$\frac{1 - |B|^2}{1 - |b|^2} = \frac{|A|^2}{|a|^2}. \quad (33.4)$$

We now substitute Eq. (33.2) to eliminate A/a from the expression above, and express B in terms of b by means of Eq. (33.3):

$$\frac{1 - |b|^2 \frac{\cos^2(i+r)}{\cos^2(i-r)}}{1 - |b|^2} = \frac{1}{\cos^2(i-r)}. \quad (33.5)$$

This equation contains only one unknown amplitude, b . Solving for $|b|^2$, we obtain

$$|b|^2 = \frac{\sin^2(i-r)}{\sin^2(i+r)} \quad (33.6)$$

and, with the aid of (33.3),

$$|B|^2 = \frac{\tan^2(i-r)}{\tan^2(i+r)}. \quad (33.7)$$

So we have found the reflection coefficient $|b|^2$ for an incident wave polarized perpendicular to the plane of incidence, and also the reflection coefficient $|B|^2$ for an incident wave polarized in the plane of incidence!

It is possible to go on with arguments of this nature and deduce that b is real. To prove this, one must consider a case where light is coming from both sides of the glass surface at the same time, a situation not easy to arrange experimentally, but fun to analyze theoretically. If we analyze this general case, we can prove that b must be real, and therefore, in fact, that $b = \pm \sin(i-r)/\sin(i+r)$. It is even possible to determine the sign by considering the case of a very, very thin layer in which there is reflection from the front and from the back surfaces, and calculating how much light is reflected. We know how much light should be reflected by a thin layer, because we know how much current is generated, and we have even worked out the fields produced by such currents.

One can show by these arguments that

$$b = -\frac{\sin(i-r)}{\sin(i+r)}, \quad B = -\frac{\tan(i-r)}{\tan(i+r)}. \quad (33.8)$$

These expressions for the reflection coefficients as a function of the angles of incidence and refraction are called Fresnel's reflection formulas.

If we consider the limit as the angles i and r go to zero, we find, for the case of normal incidence, that $B^2 \approx b^2 \approx (i-r)^2/(i+r)^2$ for both polarizations,

since the sines are practically equal to the angles, as are also the tangents. But we know that $\sin i / \sin r = n$, and when the angles are small, $i/r \approx n$. It is thus easy to show that the coefficient of reflection for normal incidence is

$$B^2 = b^2 = \frac{(n - 1)^2}{(n + 1)^2}.$$

It is interesting to find out how much light is reflected at normal incidence from the surface of water, for example. For water, n is $4/3$, so that the reflection coefficient is $(1/7)^2 \approx 2\%$. At normal incidence, only two percent of the light is reflected from the surface of water.

33-7 Anomalous refraction

The last polarization effect we shall consider was actually one of the first to be discovered: anomalous refraction. Sailors visiting Iceland brought back to Europe crystals of Iceland spar (CaCO_3) which had the amusing property of making anything seen through the crystal appear doubled, i.e., as two images. This came to the attention of Huygens, and played an important role in the discovery of polarization. As is often the case, the phenomena which are discovered first are the hardest, ultimately, to explain. It is only after we understand a physical concept thoroughly that we can carefully select those phenomena which most clearly and simply demonstrate the concept.

Anomalous refraction is a particular case of the same birefringence that we considered earlier. Anomalous refraction comes about when the optic axis, the long axis of our asymmetric molecules, is *not* parallel to the surface of the crystal. In Fig. 33-7 are drawn two pieces of birefringent material, with the optic axis as shown. In the upper figure, the incident beam falling on the material is linearly polarized in a direction perpendicular to the optic axis of the material. When this beam strikes the surface of the material, each point on the surface acts as a source of a wave which travels into the crystal with velocity v_{\perp} , the velocity of light in the crystal when the plane of polarization is normal to the optic axis. The wave-front is just the envelope or locus of all these little spherical waves, and this wave-front moves straight through the crystal and out the other side. This is just the ordinary behavior we would expect, and this ray is called the *ordinary ray*.

In the lower figure the linearly polarized light falling on the crystal has its direction of polarization turned through 90° , so that the optic axis lies in the plane of polarization. When we now consider the little waves originating at any point on the surface of the crystal, we see that they do not spread out as spherical waves. Light travelling along the optic axis travels with velocity v_{\perp} because the polarization is perpendicular to the optic axis, whereas the light travelling perpendicular to the optic axis travels with velocity v_{\parallel} because the polarization is parallel to the optic axis. In a birefringent material $v_{\parallel} \neq v_{\perp}$, and in the figure $v_{\parallel} < v_{\perp}$. A more complete analysis will show that the waves spread out on the surface of an ellipsoid, with the optic axis as major axis of the ellipsoid. The envelope of all these elliptical waves is the wavefront which proceeds through the crystal in the direction shown. Again, at the back surface the beam will be deflected just as it was at the front surface, so that the light emerges parallel to the incident beam, but displaced from it. Clearly, this beam does not follow Snell's law, but goes in an extraordinary direction. It is therefore called the *extraordinary ray*.

When an unpolarized beam strikes an anomalously refracting crystal, it is separated into an ordinary ray, which travels straight through in the normal manner, and an extraordinary ray which is displaced as it passes through the crystal. These two emergent rays are linearly polarized at right angles to each other. That this is true can be readily demonstrated with a sheet of polaroid to analyze the polarization of the emergent rays. We can also demonstrate that our interpretation of this phenomenon is correct by sending linearly polarized light into the crystal. By properly orienting the direction of polarization of the incident beam, we can make this light go straight through without splitting, or we can make it go through without splitting but with a displacement.

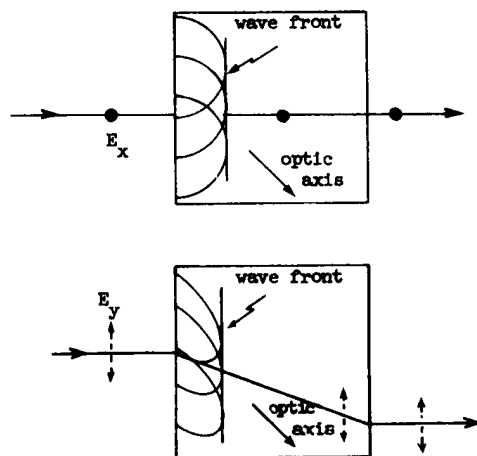


Fig. 33-7. The upper diagram shows the path of the ordinary ray through a doubly refracting crystal. The extraordinary ray is shown in the lower diagram. The optic axis lies in the plane of the paper.

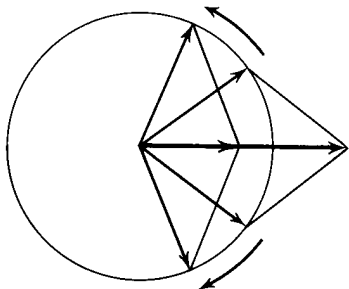


Fig. 33-8. Two oppositely rotating vectors of equal amplitude add to produce a vector in a fixed direction, but with an oscillating amplitude.

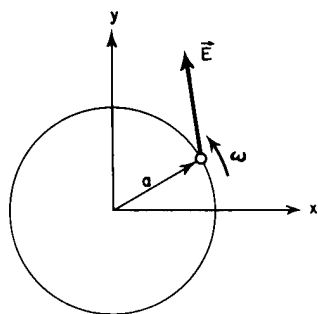


Fig. 33-9. A charge moving in a circle in response to circularly polarized light.

We have represented all the various polarization cases in Figs. 33-1 and 33-2 as superpositions of two special polarization cases, namely x and y in various amounts and phases. Other pairs could equally well have been used. Polarization along any two perpendicular axes x' , y' inclined to x and y would serve as well [for example, any polarization can be made up of superpositions of cases (a) and (e) of Fig. 33-2]. It is interesting, however, that this idea can be extended to other cases also. For example, any *linear* polarization can be made up by superposing suitable amounts at suitable phases of right and left *circular* polarizations [cases (c) and (g) of Fig. 33-2], since two equal vectors rotating in opposite directions add to give a single vector oscillating in a straight line (Fig. 33-8). If the phase of one is shifted relative to the other, the line is inclined. Thus all the pictures of Fig. 33-1 could be labeled "the superposition of equal amounts of right and left circularly polarized light at various relative phases." As the left slips behind the right in phase, the direction of the linear polarization changes. Therefore optically active materials are, in a sense, birefringent. Their properties can be described by saying that they have different indexes for right- and left-hand circularly polarized light. Superposition of right and left circularly polarized light of different intensities produces elliptically polarized light.

Circularly polarized light has another interesting property—it carries *angular momentum* (about the direction of propagation). To illustrate this, suppose that such light falls on an atom represented by a harmonic oscillator that can be displaced equally well in any direction in the plane xy . Then the x -displacement of the electron will respond to the E_x component of the field, while the y -component responds, equally, to the equal E_y component of the field but 90° behind in phase. That is, the responding electron goes around in a circle, with angular velocity ω , in response to the rotating electric field of the light (Fig. 33-9). Depending on the damping characteristics of the response of the oscillator, the direction of the displacement \vec{a} of the electron, and the direction of the force $q_e\vec{E}$ on it need not be the same but they rotate around together. The \vec{E} may have a component at right angles to \vec{a} , so work is done on the system and a torque τ is exerted. The work done per second is $\tau\omega$. Over a period of time T the energy absorbed is $\tau\omega T$, while τT is the angular momentum delivered to the matter absorbing the energy. We see therefore that a beam of right circularly polarized light containing a total energy \mathcal{E} carries an angular momentum (with vector directed along the direction of propagation) \mathcal{E}/ω . For when this beam is absorbed that angular momentum is delivered to the absorber. Left-hand circular light carries angular momentum of the opposite sign, $-\mathcal{E}/\omega$.

Relativistic Effects in Radiation

34-1 Moving sources

In the present chapter we shall describe a number of miscellaneous effects in connection with radiation, and then we shall be finished with the classical theory of light propagation. In our analysis of light, we have gone rather far and into considerable detail. The only phenomena of any consequence associated with electromagnetic radiation that we have not discussed is what happens if radiowaves are contained in a box with reflecting walls, the size of the box being comparable to a wavelength, or are transmitted down a long tube. The phenomena of so-called *cavity resonators* and *waveguides* we shall discuss later; we shall first use another physical example—sound—and then we shall return to this subject. Except for this, the present chapter is our last consideration of the classical theory of light.

We can summarize all the effects that we shall now discuss by remarking that they have to do with the effects of *moving sources*. We no longer assume that the source is localized, with all its motion being at a relatively low speed near a fixed point.

We recall that the fundamental laws of electrodynamics say that, at large distances from a moving charge, the electric field is given by the formula

$$\mathbf{E} = -\frac{q}{4\pi\epsilon_0 c^2} \frac{d^2 \mathbf{e}_{R'}}{dt^2}. \quad (34.1)$$

The second derivative of the unit vector $\mathbf{e}_{R'}$ which points in the apparent direction of the charge, is the determining feature of the electric field. This unit vector does not point toward the *present* position of the charge, of course, but rather in the direction that the charge would seem to be, if the information travels only at the finite speed c from the charge to the observer.

Associated with the electric field is a magnetic field, always at right angles to the electric field and at right angles to the apparent direction of the source, given by the formula

$$\mathbf{B} = -\mathbf{e}_{R'} \times \mathbf{E}/c. \quad (34.2)$$

Until now we have considered only the case in which motions are nonrelativistic in speed, so that there is no appreciable motion in the direction of the source to be considered. Now we shall be more general and study the case where the motion is at an arbitrary velocity, and see what different effects may be expected in those circumstances. We shall let the motion be at an arbitrary speed, but of course we shall still assume that the detector is very far from the source.

We already know from our discussion in Chapter 28 that the only things that count in $d^2 \mathbf{e}_{R'}/dt^2$ are the changes in the *direction* of $\mathbf{e}_{R'}$. Let the coordinates of the charge be (x, y, z) , with z measured along the direction of observation (Fig. 34-1). At a given moment in time, say the moment τ , the three components of the position are $x(\tau)$, $y(\tau)$, and $z(\tau)$. The distance R is very nearly equal to $R(\tau) = R_0 + z(\tau)$. Now the direction of the vector $\mathbf{e}_{R'}$ depends mainly on x and y , but hardly at all upon z : the transverse components of the unit vector are x/R and y/R , and when we differentiate these components we get things like R^2 in the denominator:

$$\frac{d(x/R)}{dt} = \frac{dx/dt}{R} - \frac{dz}{dt} \frac{x}{R^2}.$$

34-1 Moving sources

34-2 Finding the “apparent” motion

34-3 Synchrotron radiation

34-4 Cosmic synchrotron radiation

34-5 Bremsstrahlung

34-6 The Doppler effect

34-7 The ω, \mathbf{k} four-vector

34-8 Aberration

34-9 The momentum of light

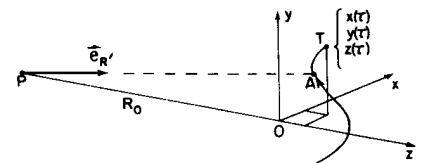


Fig. 34-1. The path of a moving charge. The true position at the time τ is at T , but the retarded position is at A .

So, when we are far enough away the only terms we have to worry about are the variations of x and y . Thus we take out the factor R_0 and get

$$\begin{aligned} E_x &= -\frac{q}{4\pi\epsilon_0 c^2 R_0} \frac{d^2 x'}{dt^2}, \\ E_y &= -\frac{q}{4\pi\epsilon_0 c^2 R_0} \frac{d^2 y'}{dt^2}, \end{aligned} \quad (34.3)$$

where R_0 is the distance, more or less, to q ; let us take it as the distance OP to the origin of the coordinates (x, y, z) . Thus the electric field is a constant multiplied by a very simple thing, the second derivatives of the x - and y -coordinates. (We could put it more mathematically by calling x and y the *transverse* components of the position vector \mathbf{r} of the charge, but this would not add to the clarity.)

Of course, we realize that the coordinates must be measured at the retarded time. Here we find that $z(\tau)$ *does* affect the retardation. What time is the retarded time? If the time of observation is called t (the time at P) then the time τ to which this corresponds at A is not the time t , but is delayed by the total distance that the light has to go, divided by the speed of light. In the first approximation, this delay is R_0/c , a constant (an uninteresting feature), but in the next approximation we must include the effects of the position in the z -direction at the time τ , because if q is a little farther back, there is a little more retardation. This is an effect that we have neglected before, and it is the only change needed in order to make our results valid for all speeds.

What we must now do is to choose a certain value of t and calculate the value of τ from it, and thus find out where x and y are at that τ . These are then the retarded x and y , which we call x' and y' , whose second derivatives determine the field. Thus τ is determined by

$$t = \tau + \frac{R_0}{c} + \frac{z(\tau)}{c}$$

and

$$x'(t) = x(\tau), \quad y'(t) = y(\tau). \quad (34.4)$$

Now these are complicated equations, but it is easy enough to make a geometrical picture to describe their solution. This picture will give us a good qualitative feeling for how things work, but it still takes a lot of detailed mathematics to deduce the precise results of a complicated problem.

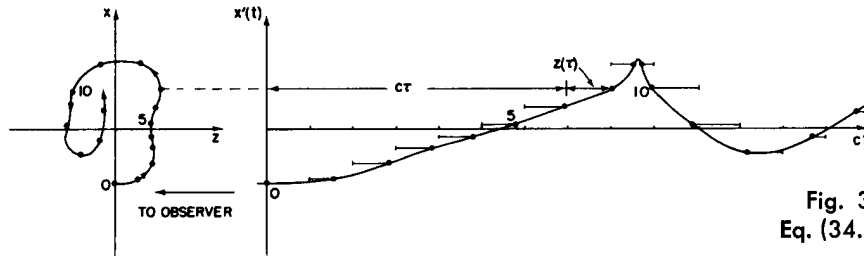


Fig. 34-2. A geometrical solution of Eq. (34.5) to find $x'(t)$.

34-2 Finding the “apparent” motion

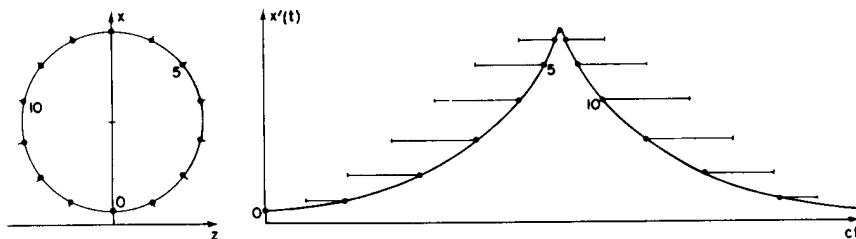
The above equation has an interesting simplification. If we disregard the uninteresting constant delay R_0/c , which just means that we must change the origin of t by a constant, then it says that

$$ct = c\tau + z(\tau), \quad x' = x(\tau), \quad y' = y(\tau). \quad (34.5)$$

Now we need to find x' and y' as functions of t , not τ , and we can do this in the following way: Eq. (34.5) says that we should take the actual motion and add a constant (the speed of light) times τ . What that turns out to mean is shown in Fig. 34-2. We take the actual motion of the charge (shown at left) and imagine that as it is going around it is being swept away from the point P at the speed c (there are no contractions from relativity or anything like that; this is just a mathematical addition of the $c\tau$). In this way we get a new motion, in which the line-

of-sight coordinate is ct , as shown at the right. (The figure shows the result for a rather complicated motion in a plane, but of course the motion may not be in one plane—it may be even more complicated than motion in a plane.) The point is that the horizontal (i.e., line-of-sight) distance now is no longer the old z , but is $z + c\tau$, and therefore is ct . Thus we have found a picture of the curve, x' (and y') against t ! All we have to do to find the field is to look at the acceleration of this curve, i.e., to differentiate it twice. So the final answer is: in order to find the electric field for a moving charge, take the motion of the charge and translate it back at the speed c to “open it out”; then the curve, so drawn, is a curve of the x' and y' positions of the function of t . The acceleration of this curve gives the electric field as a function of t . Or, if we wish, we can now imagine that this whole “rigid” curve moves forward at the speed c through the plane of sight, so that the point of intersection with the plane of sight has the coordinates x' and y' . The acceleration of this point makes the electric field. This solution is just as exact as the formula we started with—it is simply a geometrical representation.

Fig. 34-3. The $x'(t)$ curve for a particle moving at constant speed $v = 0.94c$, a circle.



If the motion is relatively slow, for instance if we have an oscillator just going up and down slowly, then when we shoot that motion away at the speed of light, we would get, of course, a simple cosine curve, and that gives a formula we have been looking at for a long time: it gives the field produced by an oscillating charge. A more interesting example is an electron moving rapidly, very nearly at the speed of light, in a circle. If we look in the plane of the circle, the retarded $x'(t)$ appears as shown in Fig. 34-3. What is this curve? If we imagine a radius vector from the center of the circle to the charge, and if we extend this radial line a little bit past the charge, just a shade if it is going fast, then we come to a point on the line that goes at the speed of light. Therefore, when we translate the motion back at the speed of light, that corresponds to having a wheel with a charge on it rolling backward (without slipping) at the speed c ; thus we find a curve which is very close to a cycloid—it is called a *hypocycloid*. If the charge is going very nearly at the speed of light, the “cusps” are very sharp indeed; if it went at exactly the speed of light, they would be actual cusps, infinitely sharp. “Infinitely sharp” is interesting; it means that near a cusp the second derivative is enormous. Once in each cycle we get a sharp pulse of electric field. This is not at all what we would get from a nonrelativistic motion, where each time the charge goes around there is an oscillation which is of about the same “strength” all the time. Instead, there are very sharp pulses of electric field spaced at time intervals $1/T_0$ apart, where T_0 is the period of revolution. These strong electric fields are emitted in a narrow cone in the direction of motion of the charge. When the charge is moving away from P , there is very little curvature and there is very little radiated field in the direction of P .

34-3 Synchrotron radiation

We have very fast electrons moving in circular paths in the synchrotron; they are travelling at very nearly the speed c , and it is possible to see the above radiation as actual *light*! Let us discuss this in more detail.

In the synchrotron we have electrons which go around in circles in a uniform magnetic field. First, let us see why they go in circles. From Eq. (12.10), we know that the force on a particle in a magnetic field is given by

$$\mathbf{F} = q\mathbf{v} \times \mathbf{B}, \quad (34.6)$$

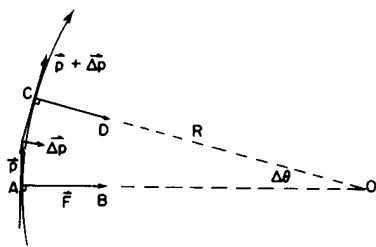


Fig. 34-4. A charged particle moves in a circular (or helical) path in a uniform magnetic field.

and it is at right angles both to the field and to the velocity. As usual, the force is equal to the rate of change of momentum with time. If the field is directed upward out of the paper, the momentum of the particle and the force on it are as shown in Fig. 34-4. Since the force is at right angles to the velocity, the kinetic energy, and therefore the speed, remains *constant*. All the magnetic field does is to change the *direction of motion*. In a short time Δt , the momentum vector changes at right angles to itself by an amount $\Delta \mathbf{p} = \mathbf{F} \Delta t$, and therefore \mathbf{p} turns through an angle $\Delta \theta = \Delta p/p = qvB \Delta t/p$, since $|\mathbf{F}| = qvB$. But in this same time the particle has gone a distance $\Delta s = v \Delta t$. Evidently, the two lines AB and CD will intersect at a point O such that $OA = OC = R$, where $\Delta s = R \Delta \theta$. Combining this with the previous expressions, we find $R \Delta \theta / \Delta t = R \omega = v = qvBR/p$, from which we find

$$p = qBR \quad (34.7)$$

and

$$\omega = qvB/p. \quad (34.8)$$

Since this same argument can be applied during the next instant, the next, and so on, we conclude that the particle must be moving in a *circle* of radius R , with angular velocity ω .

The result that the momentum of the particle is equal to a charge times the radius times the magnetic field is a very important law that is used a great deal. It is important for practical purposes because if we have elementary particles which all have the same charge and we observe them in a magnetic field, we can measure the radii of curvature of their orbits and, knowing the magnetic field, thus determine the momenta of the particles. If we multiply both sides of Eq. (34.7) by c , and express q in terms of the electronic charge, we can measure the momentum in units of the *electron volt*. In those units our formula is

$$pc(\text{ev}) = 3 \times 10^8 (q/q_e) BR, \quad (34.9)$$

where B , R , and the speed of light are all expressed in the mks system, the latter being 3×10^8 , numerically.

The mks unit of magnetic field is called a *weber per square meter*. There is an older unit which is still in common use, called a *gauss*. One weber/m^2 is equal to 10^4 gauss. To give an idea of how big magnetic fields are, the strongest magnetic field that one can usually make in iron is about 1.5×10^4 gauss; beyond that, the advantage of using iron disappears. Today, electromagnets wound with superconducting wire are able to produce steady fields of over 10^5 gauss strength—that is, 10 mks units. The field of the earth is a few tenths of a gauss at the equator.

Returning to Eq. (34.9), we could imagine the synchrotron running at a billion electron volts, so pc would be 10^9 for a billion electron volts. (We shall come back to the energy in just a moment.) Then, if we had a B corresponding to, say, 10,000 gauss, which is a good substantial field, one mks unit, then we see that R would have to be 3.3 meters. The actual radius of the Caltech synchrotron is 3.7 meters, the field is a little bigger, and the energy is 1.5 billion, but it is the same idea. So now we have a feeling for why the synchrotron has the size it has.

We have calculated the momentum, but we know that the total energy, including the rest energy, is given by $W = \sqrt{p^2 c^2 + m^2 c^4}$, and for an electron the rest energy corresponding to mc^2 is 0.511×10^6 ev, so when pc is 10^9 ev we can neglect mc^2 , and so for all practical purposes $W = pc$ when the speeds are relativistic. It is practically the same to say the energy of an electron is a billion electron volts as to say the momentum times c is a billion electron volts. If $W = 10^9$ ev, it is easy to show that the speed differs from the speed of light by but one part in eight million!

We turn now to the radiation emitted by such a particle. A particle moving on a circle of radius 3.3 meters, or 20 meters circumference, goes around once in roughly the time it takes light to go 20 meters. So the wavelength that should be emitted by such a particle would be 20 meters—in the shortwave radio region. But because of the piling up effect that we have been discussing (Fig. 34-3), and because the distance by which we must extend the radius to reach the speed c is

only one part in eight million of the radius, the cusps of the hypocycloid are enormously sharp compared with the distance between them. The acceleration, which involves a second derivative with respect to time, gets twice the “compression factor” of 8×10^6 because the time scale is reduced by eight million twice in the neighborhood of the cusp. Thus we might expect the effective wavelength to be much shorter, to the extent of 64 times 10^{12} smaller than 20 meters, and that corresponds to the x-ray region. (Actually, the cusp itself is not the entire determining factor; one must also include a certain region about the cusp. This changes the factor to the $3/2$ power instead of the square, but still leaves us above the optical region.) Thus, even though a slowly moving electron would have radiated 20-meter radiowaves, the relativistic effect cuts down the wavelength so much that we can see it! Clearly, the light should be *polarized*, with the electric field perpendicular to the uniform magnetic field.

To further appreciate what we would observe, suppose that we were to take such light (to simplify things, because these pulses are so far apart in time, we shall just take one pulse) and direct it onto a diffraction grating, which is a lot of scattering wires. After this pulse comes away from the grating, what do we see? (We should see red light, blue light, and so on, if we see any light at all.) What *do* we see? The pulse strikes the grating head-on, and all the oscillators in the grating, together, are violently moved up and then back down again, just once. They then produce effects in various directions, as shown in Fig. 34-5. But the point P is closer to one end of the grating than to the other, so at this point the electric field arrives first from wire A , next from B , and so on; finally, the pulse from the last wire arrives. In short, the sum of the reflections from all the successive wires is as shown in Fig. 34-6(a); it is an electric field which is a series of pulses, and it is very like a sine wave whose wavelength is the distance between the pulses, just as it would be for monochromatic light striking the grating! So, we get colored light all right. But, by the same argument, will we not get light from any kind of a “pulse”? No. Suppose that the curve were much smoother; then we would add all the scattered waves together, separated by a small time between them (Fig. 34-6b). Then we see that the field would not shake at all, it would be a very smooth curve, because each pulse does not vary much in the time interval between pulses.

The electromagnetic radiation emitted by relativistic charged particles circulating in a magnetic field is called *synchrotron radiation*. It is so named for obvious reasons, but it is not limited specifically to synchrotrons, or even to earthbound laboratories. It is exciting and interesting that it also occurs in nature!

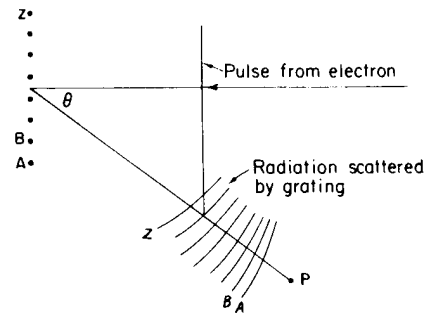


Fig. 34-5. The light which strikes a grating as a single, sharp pulse is scattered in various directions as different colors.

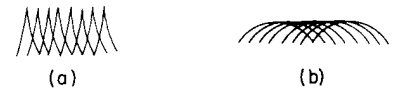


Fig. 34-6. The total electric field due to a series of (a) sharp pulses and (b) smooth pulses.

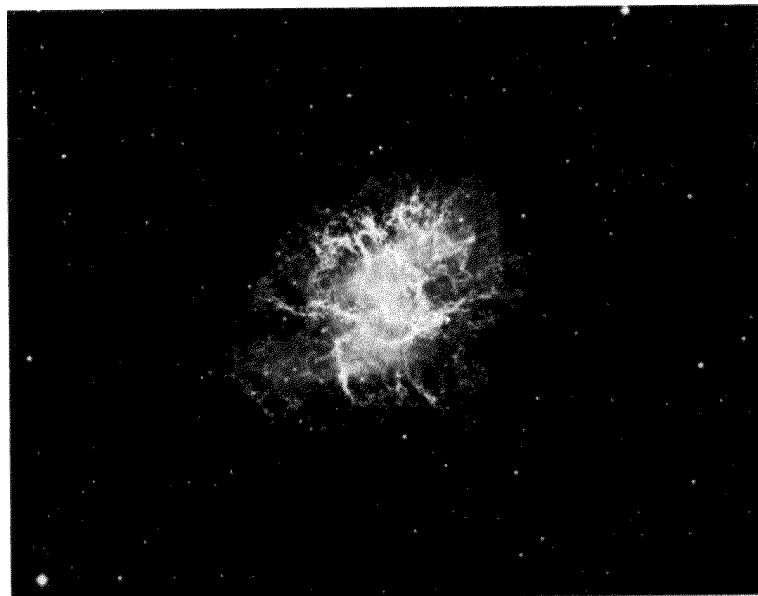


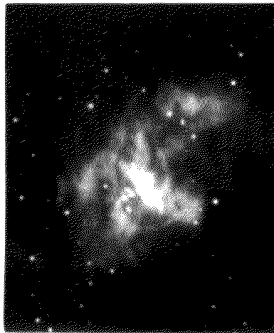
Fig. 34-7. The crab nebula as seen in all colors (no filter).

34-4 Cosmic synchrotron radiation

In the year 1054 the Chinese and Japanese civilizations were among the most advanced in the world; they were conscious of the external universe, and they recorded, most remarkably, an explosive bright star in that year. (It is amazing that none of the European monks, writing all the books of the middle ages, even bothered to write that a star exploded in the sky, but they did not.) Today we may take a picture of that star, and what we see is shown in Fig. 34-7. On the outside is a big mass of red filaments, which is produced by the atoms of the thin gas “ringing” at their natural frequencies; this makes a bright line spectrum with different frequencies in it. The red happens in this case to be due to nitrogen. On the other hand, in the central region is a mysterious, fuzzy patch of light in a *continuous* distribution of frequency, i.e., there are no special frequencies associated with particular atoms. Yet this is not dust “lit up” by nearby stars, which is one way by which one can get a continuous spectrum. We can see stars through it, so it is transparent, but it is *emitting* light.



(a)



(b)

Fig. 34-8. The crab nebula as seen through a blue filter and a polaroid. (a) Electric vector vertical. (b) Electric vector horizontal.

In Fig. 34-8 we look at the same object, using light in a region of the spectrum which has no bright spectral line, so that we see only the central region. But in this case, also, polarizers have been put on the telescope, and the two views correspond to two orientations 90° apart. We see that the pictures are different! That is to say, the light is polarized. The reason, presumably, is that there is a local magnetic field, and many very energetic electrons are going around in that magnetic field.

We have just illustrated how the electrons could go around the field in a circle. We can add to this, of course, any uniform motion in the direction of the field, since the force, $q\mathbf{v} \times \mathbf{B}$, has no component in this direction and, as we have already remarked, the synchrotron radiation is evidently polarized in a direction at right angles to the projection of the magnetic field onto the plane of sight.

Putting these two facts together, we see that in a region where one picture is bright and the other one is black, the light must have its electric field completely polarized in one direction. This means that there is a magnetic field at right angles to this direction, while in other regions, where there is a strong emission in the other picture, the magnetic field must be the other way. If we look carefully at Fig. 34-8, we may notice that there is, roughly speaking, a general set of “lines” that go one way in one picture and at right angles to this in the other. The pictures show a kind of fibrous structure. Presumably, the magnetic field lines will tend to extend relatively long distances in their own direction, and so, presumably, there are long regions of magnetic field with all the electrons spiralling one way, while in another region the field is the other way and the electrons are also spiralling that way.

What keeps the electron energy so high for so long a time? After all, it is 900 years since the explosion—how can they keep going so fast? How they maintain their energy and how this whole thing keeps going is still not thoroughly understood.

34-5 Bremsstrahlung

We shall next remark briefly on one other interesting effect of a very fast-moving particle that radiates energy. The idea is very similar to the one we have just discussed. Suppose that there are charged particles in a piece of matter and a very fast electron, say, comes by (Fig. 34-9). Then, because of the electric field around the atomic nucleus the electron is pulled, accelerated, so that the curve of its motion has a slight kink or bend in it. If the electron is travelling at very nearly the speed of light, what is the electric field produced in the direction c ? Remember our rule: we take the actual motion, translate it backwards at speed c , and that gives us a curve whose curvature measures the electric field. It was coming toward us at the speed v , so we get a backward motion, with the whole picture compressed into a smaller distance in proportion as $c - v$ is smaller than c . So, if $1 - v/c \ll 1$, there is a very sharp and rapid curvature at B' , and when we take the second derivative of that we get a very high field in the direction of the motion. So when very energetic electrons move through matter they spit radiation in a forward direction. This is called *bremsstrahlung*. As a matter of fact, the synchrotron

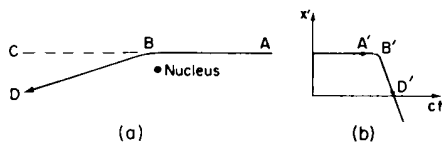


Fig. 34-9. A fast electron passing near a nucleus radiates energy in the direction of its motion.

is used, not so much to make high-energy electrons (actually if we could get them out of the machine more conveniently we would not say this) as to make very energetic photons—gamma rays—by passing the energetic electrons through a solid tungsten “target,” and letting them radiate photons from this bremsstrahlung effect.

34-6 The Doppler effect

Now we go on to consider some other examples of the effects of moving sources. Let us suppose that the source is a stationary atom which is oscillating at one of its natural frequencies, ω_0 . Then we know that the frequency of the light we would observe is ω_0 . But now let us take another example, in which we have a similar oscillator oscillating with a frequency ω_1 , and at the same time the whole atom, the whole oscillator, is moving along in a direction toward the observer at velocity v . Then the actual motion in space, of course, is as shown in Fig. 34-10(a). Now we play our usual game, we add $c\tau$; that is to say, we translate the whole curve backward and we find then that it oscillates as in Fig. 34-10(b). In a given amount of time τ , when the oscillator would have gone a distance $v\tau$, on the x' vs. ct diagram it goes a distance $(c - v)\tau$. So all the oscillations of frequency ω_1 in the time $\Delta\tau$ are now found in the interval $\Delta\tau = (1 - v/c)\Delta\tau$; they are squashed together, and as this curve comes by us at speed c , we will see light of a *higher frequency*, higher by just the compression factor $(1 - v/c)$. Thus we observe

$$\omega = \frac{\omega_1}{1 - v/c}. \quad (34.10)$$

We can, of course, analyze this situation in various other ways. Suppose that the atom were emitting, instead of sine waves, a series of pulses, pip, pip, pip, pip, at a certain frequency ω_1 . At what frequency would they be received by us? The first one that arrives has a certain delay, but the next one is delayed less because in the meantime the atom moves closer to the receiver. Therefore, the time between the “pips” is decreased by the motion. If we analyze the geometry of the situation, we find that the frequency of the pips is increased by the factor $1/(1 - v/c)$.

Is $\omega = \omega_0/(1 - v/c)$, then, the frequency that would be observed if we took an ordinary atom, which had a natural frequency ω_0 , and moved it toward the receiver at speed v ? No; as we well know, the natural frequency ω_1 of a moving atom is not the same as that measured when it is standing still, because of the relativistic dilation in the rate of passage of time. Thus if ω_0 were the true natural frequency, then the modified natural frequency ω_1 would be

$$\omega_1 = \omega_0 \sqrt{1 - v^2/c^2}. \quad (34.11)$$

Therefore the observed frequency ω is

$$\omega = \frac{\omega_0 \sqrt{1 - v^2/c^2}}{1 - v/c}. \quad (34.12)$$

The shift in frequency observed in the above situation is called the *Doppler effect*: if something moves toward us the light it emits appears more violet, and if it moves away it appears more red.

We shall now give two more derivations of this same interesting and important result. Suppose, now, that the *source* is standing still and is emitting waves at frequency ω_0 , while the *observer* is moving with speed v toward the source. After a certain period of time t the observer will have moved to a new position, a distance vt from where he was at $t = 0$. How many radians of phase will he have seen go by? A certain number, $\omega_0 t$, went past any fixed point, and in addition the observer has swept past some more by his own motion, namely a number $vt k_0$ (the number of radians per meter times the distance). So the total number of radians in the time t , or the observed frequency, would be $\omega_1 = \omega_0 + k_0 v$. We have made this analysis from the point of view of a man at rest; we would like to know how it would look to the man who is moving. Here we have to worry again about the difference in clock rate for the two observers, and this time that means that we have to *divide* by $\sqrt{1 - v^2/c^2}$. So if k_0 is the wave number, the number of radians

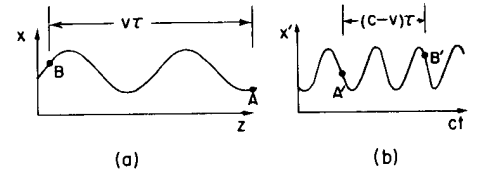


Fig. 34-10. The $x - z$ and $x' - t$ curves of a moving oscillator.

per meter in the direction of motion, and ω_0 is the frequency, then the observed frequency for a moving man is

$$\omega = \frac{\omega_0 + k_0 v}{\sqrt{1 - v^2/c^2}}. \quad (34.13)$$

For the case of light, we know that $k_0 = \omega_0/c$. So, in this particular problem, the equation would read

$$\omega = \frac{\omega_0(1 + v/c)}{\sqrt{1 - v^2/c^2}}, \quad (34.14)$$

which looks completely unlike formula (34.12)! Is the frequency that we would observe if we move toward a source different than the frequency that we would see if the source moved toward us? Of course not! The theory of relativity says that these two must be *exactly equal*. If we were expert enough mathematicians we would probably recognize that these two mathematical expressions *are* exactly equal! In fact, the *necessary* equality of the two expressions is one of the ways by which some people like to demonstrate that relativity requires a time dilation, because if we did not put those square-root factors in, they would no longer be equal.

Since we know about relativity, let us analyze it in still a third way, which may appear a little more general. (It is really the same thing, since it makes no difference *how* we do it!) According to the relativity theory there is a relationship between position and time as observed by one man and position and time as seen by another who is moving relative to him. We wrote down those relationships long ago (Chapter 16). This is the *Lorentz transformation* and its inverse:

$$\begin{aligned} x' &= \frac{x + vt}{\sqrt{1 - v^2/c^2}}, & x &= \frac{x' - vt'}{\sqrt{1 - v^2/c^2}}, \\ t' &= \frac{t + vx/c^2}{\sqrt{1 - v^2/c^2}}, & t &= \frac{t' - vx'/c^2}{\sqrt{1 - v^2/c^2}}. \end{aligned} \quad (34.15)$$

If we were standing still on the ground, the form of a wave would be $\cos(\omega t - kx)$; all the nodes and maxima and minima would follow this form. But what would a man in motion, observing the same physical wave, see? Where the field is zero, the positions of all the nodes are the same (when the field is zero, *everyone* measures the field as zero); that is a relativistic invariant. So the form is the same for the other man too, except that we must transform it into his frame of reference:

$$\cos(\omega t - kx) = \cos \left[\omega \frac{t' - vx'/c^2}{\sqrt{1 - v^2/c^2}} - k \frac{x' - vt'}{\sqrt{1 - v^2/c^2}} \right].$$

If we regroup the terms inside the brackets, we get

$$\begin{aligned} \cos(\omega t - kx) &= \cos \left[\underbrace{\frac{\omega + kv}{\sqrt{1 - v^2/c^2}}}_{\omega'} t' - \underbrace{\frac{k + v\omega/c^2}{\sqrt{1 - v^2/c^2}}}_{k'} x' \right] \\ &= \cos [\omega' t' - k' x']. \end{aligned} \quad (34.16)$$

This is again a wave, a cosine wave, in which there is a certain frequency ω' , a constant multiplying t' , and some other constant, k' , multiplying x' . We call k' the wave number, or the number of waves per meter, for the other man. Therefore the other man will see a new frequency and a new wave number given by

$$\omega' = \frac{\omega + kv}{\sqrt{1 - v^2/c^2}}, \quad (34.17)$$

$$k' = \frac{k + v\omega/c^2}{\sqrt{1 - v^2/c^2}}. \quad (34.18)$$

If we look at (34.17), we see that it is the same formula (34.13), that we obtained by a more physical argument.

34-7 The ω, k four-vector

The relationships indicated in Eqs. (34.17) and (34.18) are very interesting, because these say that the new frequency ω' is a combination of the old frequency ω and the old wave number k , and that the new wave number is a combination of the old wave number and frequency. Now the wave number is the rate of change of phase with distance, and the frequency is the rate of change of phase with time, and in these expressions we see a close analogy with the Lorentz transformation of the position and time: if ω is thought of as being like t , and k is thought of as being like x divided by c^2 , then the new ω' will be like t' , and the new k' will be like x'/c^2 . That is to say, *under the Lorentz transformation ω and k transform the same way as do t and x* . They constitute what we call a *four-vector*; when a quantity has four components transforming like time and space, it is a four-vector. Everything seems all right, then, except for one little thing: we said that a four-vector has to have *four components*; where are the other two components? We have seen that ω and k are like time and space in one space direction, but not in all directions, and so we must next study the problem of the propagation of light in three space dimensions, not just in one direction, as we have been doing up until now.

Suppose that we have a coordinate system, x, y, z , and a wave which is travelling along and whose wavefronts are as shown in Fig. 34-11. The wavelength of the wave is λ , but the direction of motion of the wave does not happen to be in the direction of one of the axes. What is the formula for such a wave? The answer is clearly $\cos(\omega t - ks)$, where $k = 2\pi/\lambda$ and s is the distance along the direction of motion of the wave—the component of the spatial position in the direction of motion. Let us put it this way: if \mathbf{r} is the vector position of a point in space, then s is $\mathbf{r} \cdot \mathbf{e}_k$, where \mathbf{e}_k is a unit vector in the direction of motion. That is, s is just $r \cos(\mathbf{r}, \mathbf{e}_k)$, the component of distance in the direction of motion. Therefore our wave is $\cos(\omega t - k\mathbf{e}_k \cdot \mathbf{r})$.

Now it turns out to be very convenient to define a vector \mathbf{k} , which is called the *wave vector*, which has a magnitude equal to the wave number, $2\pi/\lambda$, and is pointed in the direction of propagation of the waves:

$$\mathbf{k} = 2\pi\mathbf{e}_k/\lambda = k\mathbf{e}_k. \quad (34.19)$$

Using this vector, our wave can be written as $\cos(\omega t - \mathbf{k} \cdot \mathbf{r})$, or as $\cos(\omega t - k_x x - k_y y - k_z z)$. What is the significance of a component of k , say k_x ? Clearly, k_x is the rate of change of phase with respect to x . Referring to Fig. 34-11, we see that the phase changes as we change x , just as if there were a wave along x , *but of a longer wavelength*. The “wavelength in the x -direction” is longer than a natural, true wavelength by the secant of the angle α between the actual direction of propagation and the x -axis:

$$\lambda_x = \lambda / \cos \alpha. \quad (34.20)$$

Therefore the rate of change of phase, which is proportional to the *reciprocal* of λ_x , is *smaller* by the factor $\cos \alpha$; that is just how k_x would vary—it would be the magnitude of \mathbf{k} , times the cosine of the angle between \mathbf{k} and the x -axis!

That, then, is the nature of the wave vector that we use to represent a wave in three dimensions. The four quantities ω, k_x, k_y, k_z transform in relativity as a four-vector, where ω corresponds to the time, and k_x, k_y, k_z correspond to the x -, y -, and z -components of the four-vector.

In our previous discussion of special relativity (Chapter 17), we learned that there are ways of making relativistic dot products with four-vectors. If we use the position vector x_μ , where μ stands for the four components (time and three space ones), and if we call the wave vector k_μ , where the index μ again has four values, time and three space ones, then the dot product of x_μ and k_μ is written $\Sigma' k_\mu x_\mu$ (see Chapter 17). This dot product is an invariant, independent of the coordinate system; what is it equal to? By the definition of this dot product in four dimensions, it is

$$\Sigma' k_\mu x_\mu = \omega t - k_x x - k_y y - k_z z. \quad (34.21)$$

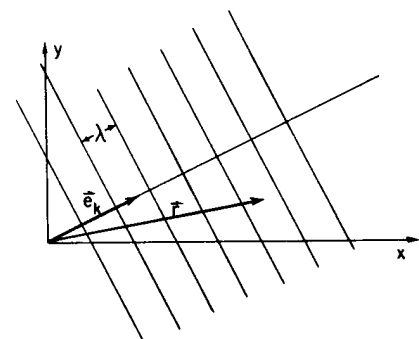


Fig. 34-11. A plane wave travelling in an oblique direction.

We know from our study of vectors that $\Sigma' k_\mu x_\mu$ is invariant under the Lorentz transformation, since k_μ is a four-vector. But this quantity is precisely what appears inside the cosine for a plane wave, and it *ought* to be invariant under a Lorentz transformation. We cannot have a formula with something that changes inside the cosine, since we know that the phase of the wave cannot change when we change the coordinate system.

34-8 Aberration

In deriving Eqs. (34.17) and (34.18), we have taken a simple example where \mathbf{k} happened to be in a direction of motion, but of course we can generalize it to other cases also. For example, suppose there is a source sending out light in a certain direction from the point of view of a man at rest, but we are moving along on the earth, say (Fig. 34-12). From which direction does the light appear to come? To find out, we will have to write down the four components of k_μ and apply the Lorentz transformation. The answer, however, can be found by the following argument: we have to point our telescope at an angle to see the light. Why? Because light is coming down at the speed c , and we are moving sidewise at the speed v , so the telescope has to be tilted forward so that as the light comes down it goes "straight" down the tube. It is very easy to see that the horizontal distance is vt when the vertical distance is ct , and therefore, if θ' is the angle of tilt, $\tan \theta' = v/c$. How nice! How nice, indeed—except for one little thing: θ' is *not* the angle at which we would have to set the telescope *relative to the earth*, because we made our analysis from the point of view of a "fixed" observer. When we said the horizontal distance is vt , the man on the earth would have found a different distance, since he measured with a "squashed" ruler. It turns out that, because of that contraction effect,

$$\tan \theta = \frac{v/c}{\sqrt{1 - v^2/c^2}}, \quad (34.22)$$

which is equivalent to

$$\sin \theta = v/c. \quad (34.23)$$

It will be instructive for the student to derive this result, using the Lorentz transformation.

This effect, that a telescope has to be tilted, is called *aberration*, and it has been observed. *How* can we observe it? Who can say where a given star *should* be? Suppose we *do* have to look in the wrong direction to see a star; how do we know it is the wrong direction? Because the earth goes around the sun. Today we have to point the telescope one way; six months later we have to tilt the telescope the other way. That is how we can tell that there is such an effect.

34-9 The momentum of light

Now we turn to a different topic. We have never, in all our discussion of the past few chapters, said anything about the effects of the *magnetic* field that is associated with light. Ordinarily, the effects of the magnetic field are very small, but there is one interesting and important effect which is a consequence of the magnetic field. Suppose that light is coming from a source and is acting on a charge and driving that charge up and down. We will suppose that the electric field is in the x -direction, so the motion of the charge is also in the x -direction: it has a position x and a velocity v , as shown in Fig. 34-13. The magnetic field is at right angles to the electric field. Now as the electric field acts on the charge and moves it up and down, what does the magnetic field do? The magnetic field acts on the charge (say an electron) only when it is moving; but the electron *is* moving, it is driven by the electric field, so the two of them work together: While the thing is going up and down it has a velocity and there is a force on it, B times v times q ; but in which *direction* is this force? *It is in the direction of the propagation of light.* Therefore, when light is shining on a charge and it is oscillating in response to that

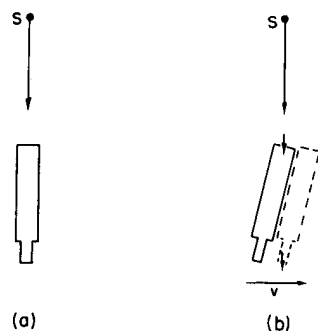


Fig. 34-12. A distant source S is viewed by (a) a stationary telescope, and (b) a laterally moving telescope.

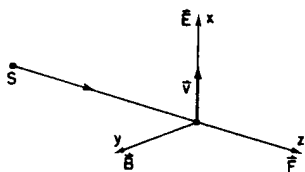


Fig. 34-13. The magnetic force on a charge which is driven by the electric field is in the direction of the light beam.

charge, there is a driving force in the direction of the light beam. This is called *radiation pressure* or light pressure.

Let us determine how strong the radiation pressure is. Evidently it is $F = qvB$ or, since everything is oscillating, it is the *time average* of this, $\langle F \rangle$. From (34.2) the strength of the magnetic field is the same as the strength of the electric field divided by c , so we need to find the average of the electric field, times the velocity, times the charge, times $1/c$: $\langle F \rangle = q\langle vE \rangle/c$. But the charge q times the field E is the electric force on a charge, and the force on the charge times the velocity is the work dW/dt being done on the charge! Therefore the force, the “pushing momentum,” that is delivered per second by the light, is equal to $1/c$ times the *energy absorbed* from the light per second! That is a general rule, since we did not say how strong the oscillator was, or whether some of the charges cancel out. *In any circumstance where light is being absorbed, there is a pressure.* The momentum that the light delivers is always equal to the energy that is absorbed, divided by c :

$$\langle F \rangle = \frac{dW/dt}{c}. \quad (34.24)$$

That light carries energy we already know. We now understand that it also carries *momentum*, and further, that the momentum carried is always $1/c$ times the energy.

When light is emitted from a source there is a recoil effect: the same thing in reverse. If an atom is emitting an energy W in some direction, then there is a recoil momentum $p = W/c$. If light is *reflected* normally from a mirror, we get twice the force.

That is as far as we shall go using the classical theory of light. Of course we know that there is a quantum theory, and that in many respects light acts like a particle. The energy of a light-particle is a constant times the frequency:

$$W = h\nu = \hbar\omega \quad (34.25)$$

We now appreciate that light also carries a momentum equal to the energy divided by c , so it is also true that these effective particles, these *photons*, carry a momentum

$$p = W/c = \hbar\omega/c = \hbar k. \quad (34.26)$$

The *direction* of the momentum is, of course, the direction of propagation of the light. So, to put it in vector form,

$$W = \hbar\omega, \quad \mathbf{p} = \hbar\mathbf{k}. \quad (34.27)$$

We also know, of course, that the energy and momentum of a particle should form a four-vector. We have just discovered that ω and \mathbf{k} form a four-vector. Therefore it is a good thing that (34.27) has the same constant in both cases; it means that the quantum theory and the theory of relativity are mutually consistent.

Equation (34.27) can be written more elegantly as $p_\mu = \hbar k_\mu$, a relativistic equation, for a particle associated with a wave. Although we have discussed this only for photons, for which k (the magnitude of \mathbf{k}) equals ω/c and $p = W/c$, the relation is much more general. In quantum mechanics *all* particles, not only photons, exhibit wavelike properties, but the frequency and wave number of the waves is related to the energy and momentum of particles by (34.27) (called the deBroglie relations) even when p is not equal to W/c .

In the last chapter we saw that a beam of right or left circularly polarized light also carries *angular momentum* in an amount proportional to the energy \mathcal{E} of the wave. In the quantum picture, a beam of circularly polarized light is regarded as a stream of photons, each carrying an angular momentum $\pm\hbar$ along the direction of propagation. That is what becomes of polarization in the corpuscular point of view—the photons carry angular momentum like spinning rifle bullets. But this “bullet” picture is really as incomplete as the “wave” picture, and we shall have to discuss these ideas more fully in a later chapter on Quantum Behavior.

Color Vision

35-1 The human eye

The phenomenon of colors depends partly on the physical world. We discuss the colors of soap films and so on as being produced by interference. But also, of course, it depends on the eye, or what happens behind the eye, in the brain. Physics characterizes the light that enters the eye, but after that, our sensations are the result of photochemical-neural processes and psychological responses.

There are many interesting phenomena associated with vision which involve a mixture of physical phenomena and physiological processes, and the full appreciation of natural phenomena, as we *see* them, must go beyond physics in the usual sense. We make no apologies for making these excursions into other fields, because the separation of fields, as we have emphasized, is merely a human convenience, and an unnatural thing. Nature is not interested in our separations, and many of the interesting phenomena bridge the gaps between fields.

In Chapter 3 we have already discussed the relation of physics to the other sciences in general terms, but now we are going to look in some detail at a specific field in which physics and other sciences are very, very closely interrelated. That area is *vision*. In particular, we shall discuss *color vision*. In the present chapter we shall discuss mainly the observable phenomena of human vision, and in the next chapter we shall consider the physiological aspects of vision, both in man and in other animals.

It all begins with the eye; so, in order to understand what phenomena we see, some knowledge of the eye is required. In the next chapter we shall discuss in some detail how the various parts of the eye work, and how they are interconnected with the nervous system. For the present, we shall describe only briefly how the eye functions (Fig. 35-1).

Light enters the eye through the *cornea*; we have already discussed how it is bent and is imaged on a layer called the *retina* in the back of the eye, so that different parts of the retina receive light from different parts of the visual field outside. The retina is not absolutely uniform: there is a place, a spot, in the center of our field of view which we use when we are trying to see things very carefully, and at which we have the greatest acuity of vision; it is called the *fovea* or *macula*. The side parts of the eye, as we can immediately appreciate from our experience in looking at things, are not as effective for seeing detail as is the center of the eye. There is also a spot in the retina where the nerves carrying all the information run out; that is a blind spot. There is no sensitive part of the retina here, and it is possible to demonstrate that if we close, say, the left eye and look straight at something, and then move a finger or another small object slowly out of the field of view it suddenly disappears somewhere. The only practical use of this fact that we know of is that some physiologist became quite a favorite in the court of a king of France by pointing this out to him; in the boring sessions that he had with his courtiers, the king could amuse himself by "cutting off their heads" by looking at one and watching another's head disappear.

Figure 35-2 shows a magnified view of the inside of the retina in somewhat schematic form. In different parts of the retina there are different kinds of structures. The objects that occur more densely near the periphery of the retina are called *rods*. Closer to the fovea, we find, besides these rod cells, also *cone* cells. We shall describe the structure of these cells later. As we get close to the fovea, the number of cones increases, and in the fovea itself there are in fact nothing but cone cells, packed very tightly, so tightly that the cone cells are much finer, or narrower

35-1 The human eye

35-2 Color depends on intensity

35-3 Measuring the color sensation

35-4 The chromaticity diagram

35-5 The mechanism of color vision

35-6 Physiochemistry of color vision

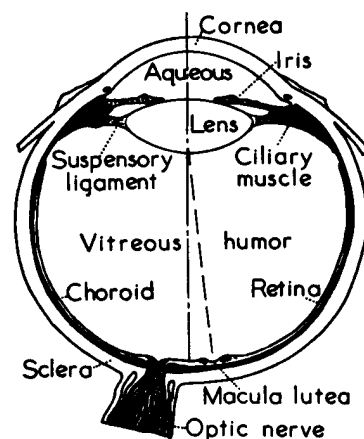


Fig. 35-1. The eye.

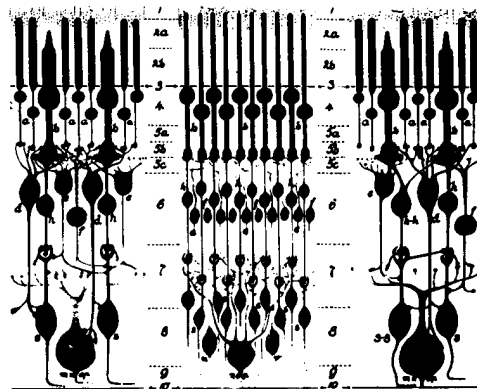


Fig. 35-2. The structure of the retina.
(Light enters from below.)

here than anywhere else. So we must appreciate that we see with the cones right in the middle of the field of view, but as we go to the periphery we have the other cells, the rods. Now the interesting thing is that in the retina each of the cells which is sensitive to light is not connected by a fiber directly to the optic nerve, but is connected to many other cells, which are themselves connected to each other. There are several kinds of cells: there are cells that carry the information toward the optic nerve, but there are others that are mainly interconnected "horizontally." There are essentially four kinds of cells, but we shall not go into these details now. The main thing we emphasize is that the light signal is already being "thought about." That is to say, the information from the various cells does not immediately go to the brain, spot for spot, but in the retina a certain amount of the information has already been digested, by a combining of the information from several visual receptors. It is important to understand that some brain-function phenomena occur in the eye itself.

35-2 Color depends on intensity

One of the most striking phenomena of vision is the dark adaptation of the eye. If we go into the dark from a brightly lighted room, we cannot see very well for a while, but gradually things become more and more apparent, and eventually we can see something where we could see nothing before. If the intensity of the light is very low, the things that we see have *no color*. It is known that this dark-adapted vision is almost entirely due to the rods, while the vision in bright light is due to the cones. As a result, there are a number of phenomena that we can easily appreciate because of this transfer of function from the cones and rods together, to just the rods.

There are many situations in which, if the light intensity were stronger, we could see color, and we would find these things quite beautiful. One example is that through a telescope we nearly always see "black and white" images of faint nebulae, but W. C. Miller of the Mt. Wilson and Palomar Observatories had the patience to make *color* pictures of some of these objects. Nobody has ever really seen these colors with the eye, but they are not artificial colors, it is merely that the light intensity is not strong enough for the cones in our eye to see them. Among the more spectacular such objects are the ring nebula and the Crab nebula. The former shows a beautiful blue inner part, with a bright red outer halo, and the latter shows a general bluish haze permeated by bright red-orange filaments.

In the bright light, apparently, the rods are at very low sensitivity but, in the dark, as time goes on they pick up their ability to see light. The variations in light intensity for which one can adapt is over a million to one. Nature does not do all this with just one kind of cell, but she passes her job from bright-light-seeing cells, the color-seeing cells, the cones, to low-intensity, dark-adapted cells, the rods. Among the interesting consequences of this shift is, first, that there is no color, and second, that there is a difference in the relative brightness of differently colored objects. It turns out that the rods see better toward the blue than the cones do, and the cones can see, for example, deep red light, while the rods find that absolutely impossible to see. So red light is black so far as the rods are concerned. Thus two pieces of colored paper, say blue and red, in which the red might be even brighter than the blue in good light, will, in the dark, appear completely reversed. It is a very striking effect. If we are in the dark and can find a magazine or something that has colors and, before we know for sure what the colors are, we judge the lighter and darker areas, and if we then carry the magazine into the light, we may see this very remarkable shift between which was the brightest color and which was not. The phenomenon is called the *Purkinje effect*.

In Fig. 35-3, the dashed curve represents the sensitivity of the eye in the dark, i.e., using the rods, while the solid curve represents it in the light. We see that the peak sensitivity of the rods is in the green region and that of the cones is more in the yellow region. If there is a red-colored page (red is about $650\text{ m}\mu$) we can see it if it is brightly lighted, but in the dark it is almost invisible.

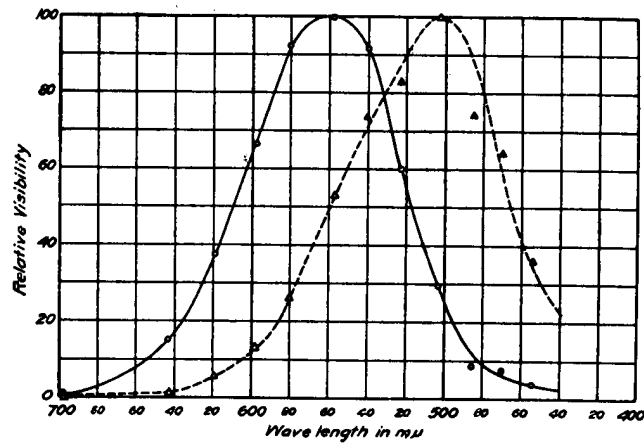


Fig. 35-3. The spectral sensitivity of the eye. Dashed curve, rods; solid curve, cones.

Another effect of the fact that rods take over in the dark, and that there are no rods in the fovea, is that when we look straight at something in the dark, our vision is not quite as acute as when we look to one side. A faint star or nebula can sometimes be seen better by looking a little to one side than directly at it, because we do not have sensitive rods in the middle of the fovea.

Another interesting effect of the fact that the number of cones decreases as we go farther to the side of the field of view is that even in a bright light color disappears as the object goes far to one side. The way to test that is to look in some particular fixed direction, let a friend walk in from one side with colored cards, and try to decide what color they are before they are right in front of you. One finds that he can see that the cards are there long before he can determine the color. When doing this, it is advisable to come in from the side opposite the blind spot, because it is otherwise rather confusing to almost see the color, then not see anything, then to see the color again.

Another interesting phenomenon is that the periphery of the retina is very sensitive to motion. Although we cannot see very well from the corner of our eye, if a little bug moves and we do not expect anything to be moving over there, we are immediately sensitive to it. We are all “wired up” to look for something jiggling to the side of the field.

35-3 Measuring the color sensation

Now we go to the cone vision, to the brighter vision, and we come to the question which is most characteristic of cone vision, and that is color. As we know, white light can be split by a prism into a whole spectrum of wavelengths which appear to us to have different colors; that is what colors are, of course: appearances. Any source of light can be analyzed by a grating or a prism, and one can determine the spectral distribution, i.e., the “amount” of each wavelength. A certain light may have a lot of blue, considerable red, very little yellow, and so on. That is all very precise in the sense of physics, but the question is, what *color* will it appear to be? It is evident that the different colors depend somehow upon the spectral distribution of the light, but the problem is to find what characteristics of the spectral distribution produce the various sensations. For example, what do we have to do to get a green color? We all know that we can simply take a piece of the spectrum which is green. But is that the *only* way to get green, or orange, or any other color?

Is there more than one spectral distribution which produces the same apparent visual effect? The answer is, definitely *yes*. There is a very limited number of visual effects, in fact just a three-dimensional manifold of them, as we shall shortly see, but there is an infinite number of different curves that we can draw for the light that comes from different sources. Now the question we have to discuss is, under what conditions do different distributions of light appear as exactly the same color to the eye?

The most powerful psycho-physical technique in color judgment is to use the eye as a *null instrument*. That is, we do not try to define what constitutes a green sensation, or to measure in what circumstances we get a green sensation, because it turns out that this is extremely complicated. Instead, we study the conditions under which two stimuli are *indistinguishable*. Then we do not have to decide whether two people see the same sensation in different circumstances, but only whether, if for one person two sensations are the same, they are also the same for another. We do not have to decide whether, when one sees something green, what it feels like inside is the same as what it feels like inside someone else when he sees something green; we do not know anything about that.

To illustrate the possibilities, we may use a series of four projector lamps which have filters on them, and whose brightnesses are continuously adjustable over a wide range: one has a red filter and makes a spot of red light on the screen, the next one has a green filter and makes a green spot, the third one has a blue filter, and the fourth one is a white circle with a black spot in the middle of it. Now if we turn on some red light, and next to it put some green, we see that in the area of overlap it produces a sensation which is not what we call reddish green, but a new color, yellow in this particular case. By changing the proportions of the red and the green, we can go through various shades of orange and so forth. If we have set it for a certain yellow, we can also obtain that same yellow, not by mixing these two colors but by mixing some other ones, perhaps a yellow filter with white light, or something like that, to get the same sensation. In other words, it is possible to make various colors in more than one way by mixing the lights from various filters.

What we have just discovered may be expressed analytically as follows. A particular yellow, for example, can be represented by a certain symbol Y , which is the "sum" of certain amounts of red-filtered light (R) and green-filtered light (G). By using two numbers, say r and g , to describe how bright the (R) and (G) are, we can write a formula for this yellow:

$$Y = rR + gG. \quad (35.1)$$

The question is, can we make *all* the different colors by adding together two or three lights of different, fixed colors? Let us see what can be done in that connection. We certainly cannot get all the different colors by mixing only red and green, because, for instance, blue never appears in such a mixture. However, by putting in some blue the central region, where all three spots overlap, may be made to appear to be a fairly nice white. By mixing the various colors and looking at this central region, we find that we can get a considerable range of colors in that region by changing the proportions, and so it is not impossible that *all* the colors can be made by mixing these three colored lights. We shall discuss to what extent this is true; it is in fact essentially correct, and we shall shortly see how to define the proposition better.

In order to illustrate our point, we move the spots on the screen so that they all fall on top of each other, and then we try to match a particular color which appears in the annular ring made by the fourth lamp. What we once thought was "white" coming from the fourth lamp now appears yellowish. We may try to match that by adjusting the red and green and blue as best we can by a kind of trial and error, and we find that we can approach rather closely this particular shade of "cream" color. So it is not hard to believe that we can make all colors. We shall try to make yellow in a moment, but before we do that, there is one color that might be very hard to make. People who give lectures on color make all the "bright" colors, but they never make *brown*, and it is hard to recall ever having seen brown light. As a matter of fact, this color is never used for any stage effect, one never sees a spotlight with brown light; so we think it might be impossible to make brown. In order to find out whether it is possible to make brown, we point out that brown light is merely something that we are not used to seeing without its background. As a matter of fact, we can make it by mixing some red and yellow. To prove that we are looking at brown light, we merely increase the brightness of the annular

background against which we see the very same light, and we see that that is, in fact, what we call brown! Brown is always a dark color next to a lighter background. We can easily change the character of the brown. For example, if we take some green out we get a reddish brown, apparently a chocolatey reddish brown, and if we put more green into it, in proportion, we get that horrible color which all the uniforms of the Army are made of, but the light from that color is not so horrible by itself; it is of yellowish green, but seen against a light background.

Now we put a yellow filter in front of the fourth light and try to match that. (The intensity must of course be within the range of the various lamps; we cannot match something which is too bright, because we do not have enough power in the lamp.) But we *can* match the yellow; we use a green and red mixture, and put in a touch of blue to make it even more perfect. Perhaps we are ready to believe that, under good conditions, we can make a perfect match of any given color.

Now let us discuss the laws of color mixture. In the first place, we found that different spectral distributions can produce the same color; next, we saw that "any" color can be made by adding together three special colors, red, blue, and green. The most interesting feature of color mixing is this: if we have a certain light, which we may call X , and if it appears indistinguishable from Y , to the eye (it may be a different spectral distribution, but it *appears* indistinguishable), we call these colors "equal," in the sense that the eye sees them as equal, and we write

$$X = Y. \quad (35.2)$$

Here is one of the great laws of color: if two spectral distributions are indistinguishable, and we *add to each one* a certain light, say Z (if we write $X + Z$, this means that we shine both lights on the same patch), and then we take Y and add the same amount of the same other light, Z , *the new mixtures are also indistinguishable*:

$$X + Z = Y + Z. \quad (35.3)$$

We have just matched our yellow; if we now shine pink light on the whole thing, it will still match. So adding any other light to the matched lights leaves a match. In other words, we can summarize all these color phenomena by saying that once we have a match between two colored lights, seen next to each other in the same circumstances, then this match will remain, and one light can be substituted for the other light in any other color mixing situation. In fact, it turns out, and it is very important and interesting, that this matching of the color of lights is not dependent upon the characteristics of the eye at the moment of observation: we know that if we look for a long time at a bright red surface, or a bright red light, and then look at a white paper, it looks greenish, and other colors are also distorted by our having looked so long at the bright red. If we now have a match between, say, two yellows, and we look at them and make them match, then we look at a bright red surface for a long time, and then turn back to the yellow, it may not look yellow any more; I do not know what color it will look, but it will not look yellow. Nevertheless *the yellows will still look matched*, and so, as the eye adapts to various levels of intensity, the color match still works, with the obvious exception of when we go into the region where the intensity of the light gets so low that we have shifted from cones to rods; then the color match is no longer a color match, because we are using a different system.

The second principle of color mixing of lights is this: *any color at all can be made from three different colors*, in our case, red, green, and blue lights. By suitably mixing the three together we can make anything at all, as we demonstrated with our two examples. Further, these laws are very interesting mathematically. For those who are interested in the mathematics of the thing, it turns out as follows. Suppose that we take our three colors, which were red, green, and blue, but label them A , B , and C , and call them our *primary* colors. Then any color could be made by certain amounts of these three: say an amount a of color A , an amount b of color B , and an amount c of color C makes X :

$$X = aA + bB + cC. \quad (35.4)$$

Now suppose another color Y is made from the same three colors:

$$Y = a'A + b'B + c'C. \quad (35.5)$$

Then it turns out that the mixture of the two lights (it is one of the consequences of the laws that we have already mentioned) is obtained by taking the sum of the components of X and Y :

$$Z = X + Y = (a + a')A + (b + b')B + (c + c')C. \quad (35.6)$$

It is just like the mathematics of the addition of vectors, where (a, b, c) are the components of one vector, and (a', b', c') are those of another vector, and the new light Z is then the “sum” of the vectors. This subject has always appealed to physicists and mathematicians. In fact, Schrödinger wrote a wonderful paper on color vision in which he developed this theory of vector analysis as applied to the mixing of colors.

Now a question is, what are the correct primary colors to use? There is no such thing as “the” correct primary colors for the mixing of lights. There may be, for practical purposes, three paints that are more useful than others for getting a greater variety of mixed pigments, but we are not discussing that matter now. *Any three differently colored lights whatsoever** can always be mixed in the correct proportion to produce *any color whatsoever*. Can we demonstrate this fantastic fact? Instead of using red, green, and blue, let us use red, blue, and yellow in our projector. Can we use red, blue, and yellow to make, say, green?

By mixing these three colors in various proportions, we get quite an array of different colors, ranging over quite a spectrum. But as a matter of fact, after a lot of trial and error, we find that nothing ever looks like green. The question is, *can* we make green? The answer is yes. How? *By projecting some red onto the green*, then we can make a match with a certain mixture of yellow and blue! So we have matched them, except that we had to cheat by putting the red on the other side. But since we have some mathematical sophistication, we can appreciate that what we really showed was not that X could always be made, say, of red, blue, and yellow, but by putting the red on the other side we found that red plus X could be made out of blue and yellow. Putting it on the other side of the equation, we can interpret that as a *negative amount*, so if we will allow that the coefficients in equations like (35.4) can be both positive and negative, and if we interpret negative amounts to mean that we have to *add* those to the *other side*, then any color can be matched by any three, and there is no such thing as “the” fundamental primaries.

We may ask whether there are three colors that come only with positive amounts for all mixings. The answer is no. Every set of three primaries requires negative amounts for some colors, and therefore there is no unique way to define a primary. In elementary books they are said to be red, green, and blue, but that is merely because with these a *wider range* of colors is available without minus signs for some of the combinations.

35-4 The chromaticity diagram

Now let us discuss the combination of colors on a mathematical level as a geometrical proposition. If any one color is represented by Eq. (35.4), we can plot it as a vector in space by plotting along three axes the amounts a , b , and c , and then a certain color is a point. If another color is a' , b' , c' , that color is located somewhere else. The sum of the two, as we know, is the color which comes from adding these as vectors. We can simplify this diagram and represent everything on a plane by the following observation: if we had a certain color light, and merely doubled a and b and c , that is, if we make them all stronger in the same ratio, it is the same color, but brighter. So if we agree to reduce everything to the *same light intensity*, then we can project everything onto a plane, and this has been done in Fig. 35-4. It follows that any color obtained by mixing a given two in some pro-

* Except, of course, if one of the three can be matched by mixing the other two.

portion will lie somewhere on a line drawn between the two points. For instance, a fifty-fifty mixture would appear halfway between them, and $1/4$ of one and $3/4$ of the other would appear $1/4$ of the way from one point to the other, and so on. If we use a blue and a green and a red, as primaries, we see that all the colors that we can make with positive coefficients are inside the dotted triangle, which contains almost all of the colors that we can ever see, because all the colors that we can ever see are enclosed in the oddly shaped area bounded by the curve. Where did this area come from? Once somebody made a very careful match of all the colors that we can see against three special ones. But we do not have to check *all* colors that we can see, we only have to check the pure spectral colors, the lines of the spectrum. Any light can be considered as a sum of various positive amounts of various pure spectral colors—pure from the physical standpoint. A given light will have a certain amount of red, yellow, blue, and so on—spectral colors. So if we know how much of each of our three chosen primaries is needed to make each of these pure components, we can calculate how much of each is needed to make our given color. So, if we find out what the *color coefficients* of all the spectral colors are for any given three primary colors, then we can work out the whole color mixing table.

An example of such experimental results for mixing three lights together is given in Fig. 35-5. This figure shows the amount of each of three different particular primaries, red, green and blue, which is required to make each of the spectral colors. Red is at the left end of the spectrum, yellow is next, and so on, all the way to blue. Notice that at some points minus signs are necessary. It is from such data that it is possible to locate the position of all of the colors on a chart, where the x - and the y -coordinates are related to the amounts of the different primaries that are used. That is the way that the curved boundary line has been found. It is the locus of the pure spectral colors. Now any other color can be made by adding spectral lines, of course, and so we find that anything that can be produced by connecting one part of this curve to another is a color that is available in nature. The straight line connects the extreme violet end of the spectrum with the extreme red end. It is the locus of the purples. Inside the boundary are colors that can be made with lights, and outside it are colors that cannot be made with lights, and nobody has ever seen them (except, possibly, in after-images!).

35-5 The mechanism of color vision

Now the next aspect of the matter is the question, *why* do colors behave in this way? The simplest theory, proposed by Young and Helmholtz, supposes that in the eye there are three different pigments which receive the light and that these have different absorption spectra, so that one pigment absorbs strongly, say, in the red, another absorbs strongly in the blue, another absorbs in the green. Then when we shine a light on them we will get different amounts of absorptions in the three regions, and these three pieces of information are somehow maneuvered in the brain or in the eye, or somewhere, to decide what the color is. It is easy to demonstrate that all of the rules of color mixing would be a consequence of this proposition. There has been considerable debate about the thing because the next problem, of course, is to find the absorption characteristics of each of the three pigments. It turns out, unfortunately, that because we can transform the color coordinates in any manner we want to, we can only find all kinds of linear combinations of absorption curves by the color-mixing experiments, but not the curves for the individual pigments. People have tried in various ways to obtain a specific curve which does describe some particular physical property of the eye. One such curve is called a *brightness curve*, demonstrated in Fig. 35-3. In this figure are two curves, one for eyes in the dark, the other for eyes in the light; the latter is the cone brightness curve. This is measured by finding what is the smallest amount of colored light we need in order to be able to just see it. This measures how sensitive the eye is in different spectral regions. There is another very interesting way to measure this. If we take two colors and make them appear in an area, by flickering back and forth from one to the other, we see a flicker if the frequency is too low. However, as the frequency increases,

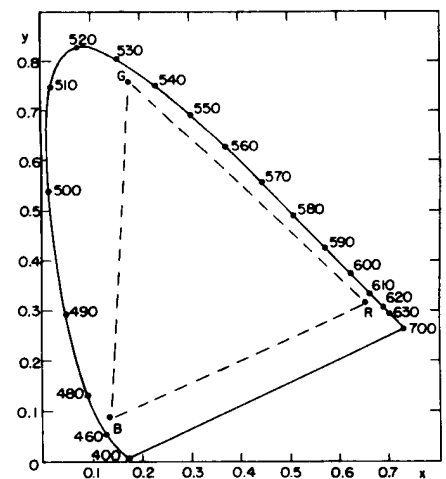


Fig. 35-4. The standard chromaticity diagram.

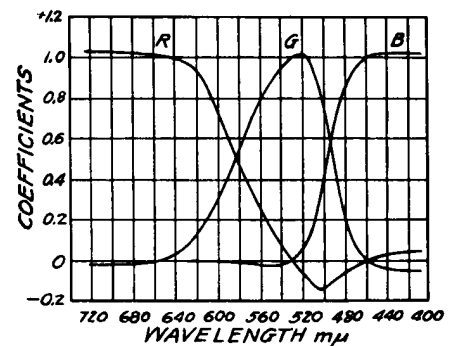


Fig. 35-5. The color coefficients of pure spectral colors in terms of a certain set of standard primary colors.

the flicker will ultimately disappear at a certain frequency that depends on the brightness of the light, let us say at 16 repetitions per second. Now if we adjust the brightness or the intensity of one color against the other, there comes an intensity where the flicker at 16 cycles disappears. To get flicker with the brightness so adjusted, we have to go to a much lower frequency in order to see a flicker of the color. So, we get what we call a flicker of the brightness at a higher frequency and, at a lower frequency, a flicker of the color. It is possible to match two colors for "equal brightness" by this flicker technique. The results are almost, but not exactly, the same as those obtained by measuring the threshold sensitivity of the eye for seeing weak light by the cones. Most workers use the flicker system as a definition of the brightness curve.

Now, if there are three color-sensitive pigments in the eye, the problem is to determine the shape of the absorption spectrum of each one. How? We know there are people who are color blind—eight percent of the male population, and one-half of one percent of the female population. Most of the people who are color blind or abnormal in color vision have a different degree of sensitivity than others to a variation of color, but they still need three colors to match. However, there are some who are called *dichromats*, for whom any color can be matched using only *two* primary colors. The obvious suggestion, then, is to say that they are missing one of the three pigments. If we can find three kinds of color-blind dichromats who have different color-mixing rules, one kind should be missing the *red*, another the *green*, and another the *blue* pigmentation. By measuring all these types we can determine the three curves! It turns out that there *are* three types of dichromatic color blindness; there are two common types and a third very rare type, and from these three it has been possible to deduce the pigment absorption spectra.

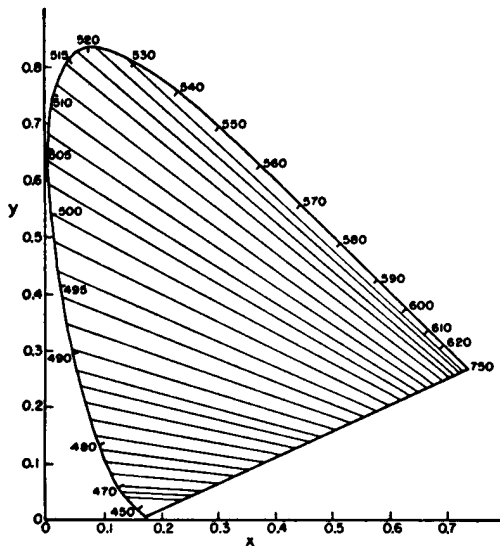


Fig. 35-6. Loci of colors confused by deuteranopes.

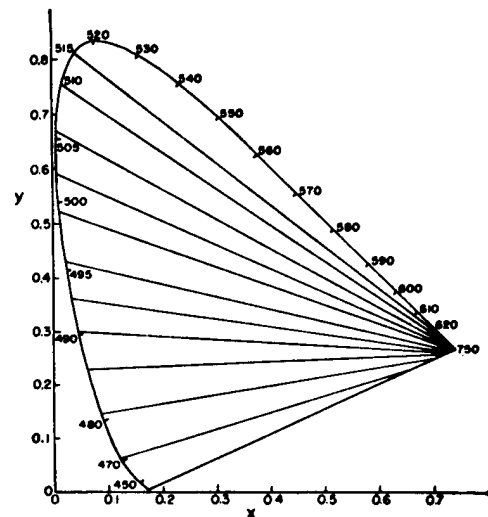


Fig. 35-7. Loci of colors confused by protanopes.

Figure 35-6 shows the color mixing of a particular type of color-blind person called a deuteranope. For him, the loci of constant colors are not points, but certain lines, along each of which the color appears to him to be the same. If the theory that he is missing one of the three pieces of information is right, all these lines should intersect at a point. If we carefully measure on this graph, they *do* intersect perfectly. Obviously, therefore, this has been made by a mathematician and does not represent real data! As a matter of fact, if we look at the latest paper with real data, it turns out that in the graph of Fig. 35-6, the point of focus of all the lines is not exactly at the right place. Using the lines in the above figure, we cannot find reasonable spectra; we need negative and positive absorptions in different regions. But using the new data of Yustova, it turns out that each of the absorption curves is everywhere positive.

Figure 35-7 shows a different kind of color blindness, that of the protanope, which has a focus near the red end of the boundary curve. Yustova gets approximately the same position in this case. Using the three different kinds of color blindness, the three pigment response curves have finally been determined, and are shown in Fig. 35-8. Finally? Perhaps. There is a question as to whether the three-pigment idea is right, whether color blindness results from lack of one pigment, and even whether the color-mix data on color blindness are right. Different workers get different results. This field is still very much under development.

35-6 Physiochemistry of color vision

Now, what about checking these curves against actual pigments in the eye? The pigments that can be obtained from a retina consist mainly of a pigment called *visual purple*. The most remarkable features of this are, first, that it is in the eye of almost every vertebrate animal, and second, that its response curve fits beautifully with the sensitivity of the eye, as seen in Fig. 35-9, in which are plotted on the same scale the absorption of visual purple and the sensitivity of the dark-adapted eye. This pigment is evidently the pigment that we see with in the dark: visual purple is the pigment for the rods, and it has nothing to do with color vision. This fact was discovered in 1877. Even today it can be said that the color pigments of the cones have never been obtained in a test tube. In 1958 it could be said that the color pigments had never been seen at all. But since that time, two of them have been detected by Rushton by a very simple and beautiful technique.

The trouble is, presumably, that since the eye is so weakly sensitive to bright light compared with light of low intensity, it needs a lot of visual purple to see with, but not much of the color pigments for seeing colors. Rushton's idea is to *leave the pigment in the eye*, and measure it anyway. What he does is this. There is an instrument called an ophthalmoscope for sending light into the eye through the lens and then focusing the light that comes back out. With it one can measure how much is reflected. So one measures the reflection coefficient of light which has gone *twice* through the pigment (reflected by a back layer in the eyeball, and coming out through the pigment of the cone again). Nature is not always so beautifully designed. The cones are interestingly designed so that the light that comes into the cone bounces around and works its way down into the little sensitive points at the apex. The light goes right down into the sensitive point, bounces at the bottom and comes back out again, having traversed a considerable amount of the color-vision pigment; also, by looking at the fovea, where there are no rods, one is not confused by visual purple. But the color of the retina has been seen a long time ago: it is a sort of orangey pink; then there are all the blood vessels, and the color of the material at the back, and so on. How do we know when we are looking at the pigment? *Answer:* First we take a color-blind person, who has fewer pigments and for whom it is therefore easier to make the analysis. Second, the various pigments, like visual purple, have an intensity change when they are bleached by light; when we shine light on them they change their concentration. So, while looking at the absorption spectrum of the eye, Rushton put *another* beam in the whole eye, which changes the concentration of the pigment, and he measured the *change* in the spectrum, and the difference, of course, has nothing to do with the amount of blood or the color of the reflecting layers, and so on, but only the pigment, and in this manner Rushton obtained a curve for the pigment of the protanope eye, which is given in Fig. 35-10.

The second curve in Fig. 35-10 is a curve obtained with a normal eye. This was obtained by taking a normal eye and, having already determined what one pigment was, bleaching the other one in the red where the first one is insensitive. Red light has no effect on the protanope eye, but does in the normal eye, and thus one can obtain the curve for the missing pigment. The shape of one curve fits beautifully with Yustova's green curve, but the red curve is a little bit displaced. So perhaps we are getting on the right track. Or perhaps not—the latest work with deuteranopes does not show any definite pigment missing.

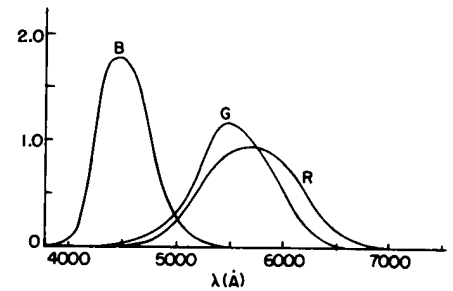


Fig. 35-8. The spectral sensitivity curves of a normal trichromat's receptors.

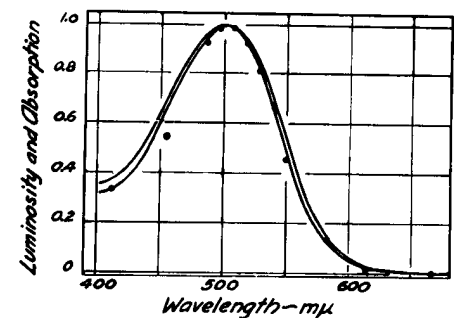


Fig. 35-9. The sensitivity curve of the dark-adapted eye, compared with the absorption curve of visual purple.

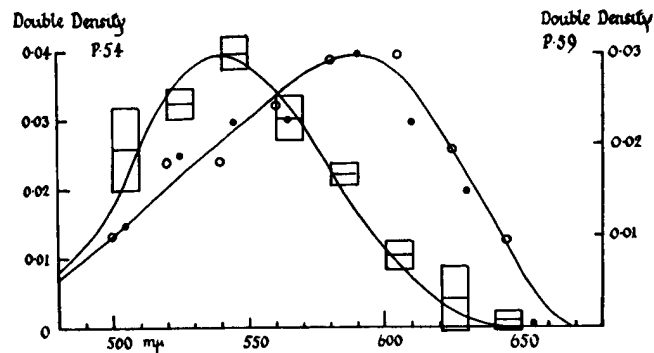


Fig. 35-10. Absorption spectrum of the color pigment of a protanope color-blind eye (squares) and a normal eye (dots).

Color is not a question of the physics of the light itself. Color is a *sensation*, and the sensation for different colors is different in different circumstances. For instance, if we have a pink light, made by superimposing crossing beams of white light and red light (all we can make with white and red is pink, obviously), we may show that white light may appear blue. If we place an object in the beams, it casts two shadows—one illuminated by the white light alone and the other by the red. For most people the “white” shadow of an object looks blue, but if we keep expanding this shadow until it covers the entire screen, we see that it suddenly appears white, not blue! We can get other effects of the same nature by mixing red, yellow, and white light. Red, yellow, and white light can produce only orangey yellows, and so on. So if we mix such lights roughly equally, we get only orange light. Nevertheless, by casting different kinds of shadows in the light, with various overlaps of colors, one gets quite a series of beautiful colors which are not in the light themselves (that is only orange), but in our *sensations*. We clearly *see* many different colors that are quite unlike the “physical” ones in the beam. It is very important to appreciate that a retina is already “thinking” about the light; it is comparing what it sees in one region with what it sees in another, although not consciously. What we know of how it does that is the subject of the next chapter.

BIBLIOGRAPHY

- Committee on Colorimetry, Optical Society of America, *The Science of Color*, Thomas Y. Crowell Company, New York, 1953.
- HECHT, S., S. SHLAER, and M. H. PIRENNE, “Energy, Quanta, and Vision,” *Journal of General Physiology*, 1942, **25**, 819–840.
- MORGAN, CLIFFORD and ELIOT STELLAR, *Physiological Psychology*, 2nd ed., McGraw-Hill Book Company, Inc., 1950.
- NUBERG, N. D. and E. N. YUSTOVA, “Researches on Dichromatic Vision and the Spectral Sensitivity of the Receptors of Trichromats,” presented at Symposium No. 8, *Visual Problems of Colour*, Vol. II, National Physical Laboratory, Teddington, England, September 1957. Published by Her Majesty’s Stationery Office, London, 1958.
- RUSHTON, W. A., “The Cone Pigments of the Human Fovea in Colour Blind and Normal,” presented at Symposium No. 8, *Visual Problems of Colour*, Vol. I, National Physical Laboratory, Teddington, England, September 1957. Published by Her Majesty’s Stationery Office, London, 1958.
- WOODWORTH, ROBERT S., *Experimental Psychology*, Henry Holt and Company, New York, 1938. Revised edition, 1954, by Robert S. Woodworth and H. Schlosberg.

Mechanisms of Seeing

36-1 The sensation of color

In discussing the sense of sight, we have to realize that (outside of a gallery of modern art!) one does not see random spots of color or spots of light. When we look at an object we see a *man* or a *thing*; in other words, the brain interprets what we see. How it does that, no one knows, and it does it, of course, at a very high level. Although we evidently do learn to recognize what a man looks like after much experience, there are a number of features of vision which are more elementary but which also involve combining information from different parts of what we see. To help us understand how we make an interpretation of an entire image, it is worth while to study the earliest stages of the putting together of information from the different retinal cells. In the present chapter we shall concentrate mainly on that aspect of vision, although we shall also mention a number of side issues as we go along.

An example of the fact that we have an accumulation, at a very elementary level, of information from several parts of the eye at the same time, beyond our voluntary control or ability to learn, was that blue shadow which was produced by white light when both white and red were shining on the same screen. This effect at least involves the knowledge that the background of the screen is pink, even though, when we are looking at the blue shadow, it is only "white" light coming into a particular spot in the eye; somewhere, pieces of information have been put together. The more complete and familiar the context is, the more the eye will make corrections for peculiarities. In fact, Land has shown that if we mix that apparent blue and the red in various proportions, by using two photographic transparencies with absorption in front of the red and the white in different proportions, it can be made to represent a real scene, with real objects, rather faithfully. In this case we get a lot of intermediate apparent colors too, analogous to what we would get by mixing red and blue-green; it seems to be an almost complete set of colors, but if we look very hard at them, they are not so very good. Even so, it is surprising how much can be obtained from just red and white. The more the scene looks like a real situation, the more one is able to compensate for the fact that all the light is actually nothing but pink!

Another example is the appearance of "colors" in a black-and-white rotating disc, whose black and white areas are as shown in Fig. 36-1. When the disc is rotated, the variations of light and dark at any one radius are exactly the same; it is only the background that is different for the two kinds of "stripes." Yet one of the "rings" appears colored with one color and the other with another.* No one yet understands the reason for those colors, but it is clear that information is being put together at a very elementary level, in the eye itself, most likely.

Almost all present-day theories of color vision agree that the color-mixing data indicate that there are only three pigments in the cones of the eye, and that it is the spectral absorption in those three pigments that fundamentally produces the color sense. But the total sensation that is associated with the absorption characteristics of the three pigments acting together is not necessarily the sum of the individual sensations. We all agree that yellow simply *does not* seem to be reddish green; in fact it might be a tremendous surprise to most people to discover that light is, in fact, a mixture of colors, because presumably the sensation of light

36-1 The sensation of color

36-2 The physiology of the eye

36-3 The rod cells

36-4 The compound (insect) eye

36-5 Other eyes

36-6 Neurology of vision

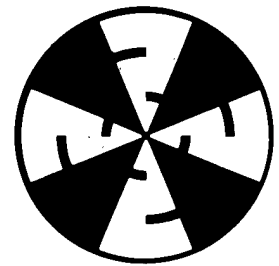


Fig. 36-1. When a disc like the above is spun, colors appear in only one of the two darker "rings." If the spin direction is reversed, the colors appear in the other ring.

* The colors depend on speed of rotation, on the brightness of illumination, and to some extent on who looks at them and how intently he stares at them.

is due to some other process than a simple mixture like a chord in music, where the three notes are there at the same time and if we listen hard we can hear them individually. We cannot look hard and see the red and the green.

The earliest theories of vision said that there are three pigments and three kinds of cones, each kind containing one pigment; that a nerve runs from each cone to the brain, so that the three pieces of information are carried to the brain; and then in the brain, anything can happen. This, of course, is an incomplete idea: it does no good to discover that the information is carried along the optic nerve to the brain, because we have not even started to solve the problem. We must ask more basic questions: Does it make any difference *where* the information is put together? Is it important that it be carried right up into the brain in the optic nerve, or could the retina do some analysis first? We have seen a picture of the retina as an extremely complicated thing with lots of interconnections (Fig. 35-2) and it might make some analyses.

As a matter of fact, people who study anatomy and the development of the eye have shown that the retina is, in fact, the brain: in the development of the embryo, a piece of the brain comes out in front, and long fibers grow back, connecting the eyes to the brain. The retina is organized in just the way the brain is organized and, as someone has beautifully put it, "The brain has developed a way to look out upon the world." The eye is a piece of brain that is touching light, so to speak, on the outside. So it is not at all unlikely that some analysis of the color has already been made in the retina.

This gives us a very interesting opportunity. None of the other senses involves such a large amount of calculation, so to speak, before the signal gets into a nerve that one can make measurements on. The calculations for all the rest of the senses usually happen in the brain itself, where it is very difficult to get at specific places to make measurements, because there are so many interconnections. Here, with the visual sense, we have the light, three layers of cells making calculations, and the results of the calculations being transmitted through the optic nerve. So we have the first chance to observe physiologically how, perhaps, the first layers of the brain work in their first steps. It is thus of double interest, not simply interesting for vision, but interesting to the whole problem of physiology.

The fact that there are three pigments does not mean that there must be three kinds of sensations. One of the other theories of color vision has it that there are really opposing color schemes (Fig. 36-2). That is, one of the nerve fibers carries a lot of impulses if there is yellow being seen, and less than usual for blue. Another nerve fiber carries green and red information in the same way, and another, white and black. In other words, in this theory someone has already started to make a guess as to the system of wiring, the method of calculation.

The problems we are trying to solve by guessing at these first calculations are questions about the apparent colors that are seen on a pink background, what happens when the eye is adapted to different colors, and also the so-called psychological phenomena. The psychological phenomena are of the nature, for instance, that white does not "feel" like red and yellow and blue, and this theory was advanced because the psychologists say that there are *four* apparent pure colors: "There are four stimuli which have a remarkable capacity to evoke psychologically simple blue, yellow, green, and red hues respectively. Unlike sienna, magenta, purple, or most of the discriminable colors, these simple hues are unmixed in the sense that none partakes of the nature of the other; specifically, blue is not yellowish, reddish, or greenish, and so on; these are psychologically primary hues." That is a psychological fact, so-called. To find out from what evidence this psychological fact was deduced, we must search very hard indeed through all the literature. In the modern literature all we find on the subject are repeats of the same statement, or of one by a German psychologist, who uses as one of his authorities Leonardo da Vinci, who, of course, we all know was a great artist. He says, "Leonardo thought there were five colors." Then, looking still further, we find, in a still older book, the evidence for the subject. The book says something like this: "Purple is reddish-blue, orange is reddish-yellow, but can red be seen as purplish-orange? Are not red and yellow more unitary than purple or orange? The average person, asked to state which

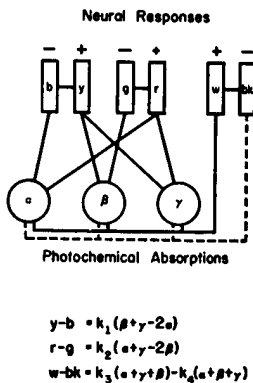


Fig. 36-2. Neural connections according to an "opponent" theory of color vision.

colors are unitary, names red, yellow, and blue, these three, and some observers add a fourth, green. Psychologists are accustomed to accept the four as salient hues." So that is the situation in the psychological analysis of this matter: if everybody says there are three, and somebody says there are four, and they want it to be four, it will be four. That shows the difficulty with psychological researches. It is clear that we have such feelings, but it is very difficult to obtain much information about them.

So the other direction to go is the physiological direction, to find out experimentally what actually happens in the brain, the eye, the retina, or wherever, and perhaps to discover that some combinations of impulses from various cells move along certain nerve fibers. Incidentally, primary pigments do not have to be in separate cells; one could have cells in which are mixtures of the various pigments, cells with the red and the green pigments, cells with all three (the information of all three is then white information), and so on. There are many ways of hooking the system up, and we have to find out which way nature has used. It would be hoped, ultimately, that when we understand the physiological connections we will have a little bit of understanding of some of those aspects of the psychology, so we look in that direction.

36-2 The physiology of the eye

We begin by talking not only about color vision, but about vision in general, just to remind ourselves about the interconnections in the retina, shown in Fig. 35-2. The retina is really like the surface of the brain. Although the actual picture through a microscope is a little more complicated looking than this somewhat schematized drawing, by careful analysis one can see all these interconnections. There is no question that one part of the surface of the retina is connected to other parts, and that the information that comes out on the long axons, which produce the optic nerve, are combinations of information from many cells. There are three layers of cells in the succession of function: there are retinal cells, which are the ones that the light affects, an intermediate cell which takes information from a single or a few retinal cells and gives it out again to several cells in a third layer of cells and carries it to the brain. There are all kinds of cross connections between cells in the layers.

We now turn to some aspects of the structure and performance of the eye (see Fig. 35-1). The focusing of the light is accomplished mainly by the cornea, by the fact that it has a curved surface which "bends" the light. This is why we cannot see clearly under water, because we then do not have enough difference between the index of the cornea, which is 1.37, and that of the water, which is 1.33. Behind the cornea is water, practically, with an index of 1.33, and behind that is a lens which has a very interesting structure: it is a series of layers, like an onion, except that it is all transparent, and it has an index of 1.40 in the middle and 1.38 at the outside. (It would be nice if we could make optical glass in which we could adjust the index throughout, for then we would not have to curve it as much as we do when we have a uniform index.) Furthermore, the shape of the cornea is not that of a sphere. A spherical lens has a certain amount of spherical aberration. The cornea is "flatter" at the outside than is a sphere, in just such a manner that the spherical aberration is less for the cornea than it would be if we put a spherical lens in there! The light is focused by the cornea-lens system onto the retina. As we look at things that are closer and farther away, the lens tightens and loosens and changes the focus to adjust for the different distances. To adjust for the total amount of light there is the iris, which is what we call the color of the eye, brown or blue, depending on who it is; as the amount of light increases and decreases, the iris moves in and out.

Let us now look at the neural machinery for controlling the accommodation of the lens, the motion of the eye, the muscles which turn the eye in the socket, and the iris, shown schematically in Fig. 36-3. Of all the information that comes out of the optic nerve *A*, the great majority is divided into one of two bundles (which we will talk about later) and thence to the brain. But there are a few fibers, of

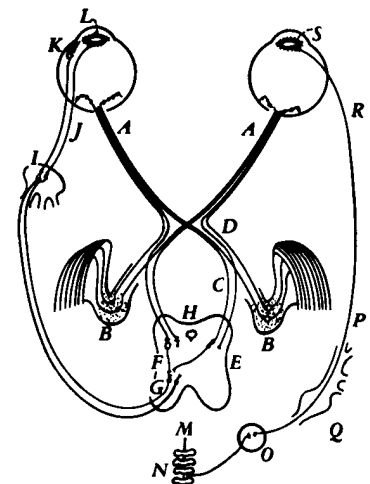


Fig. 36-3. The neural interconnections for the mechanical operation of the eyes.

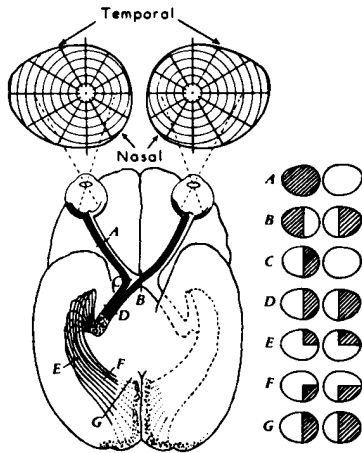


Fig. 36-4. The neural connections from the eyes to the visual cortex.

interest to us now, which do not run directly to the visual cortex of the brain where we "see" the images, but instead go into the mid-brain *H*. These are the fibers which measure the average light and make adjustment for the iris; or, if the image looks foggy, they try to correct the lens; or, if there is a double image, they try to adjust the eye for binocular vision. At any rate, they go through the mid-brain and feed back into the eye. At *K* are the muscles which run the accommodation of the lens, and at *L* another one that runs into the iris. The iris has two muscle systems. One is a circular muscle *L* which, when it is excited, pulls in and closes down the iris; it acts very rapidly and the nerves are directly connected from the brain through short axons into the iris. The opposite muscles are radial muscles, so that, when the things get dark and the circular muscle relaxes, these radial muscles pull out. Here we have, as in many places in the body, a pair of muscles which work in opposite directions, and in almost every such case the nerve systems which control the two are very delicately adjusted, so that when signals are sent in to tighten one, signals are automatically sent in to loosen the other. The iris is a peculiar exception: the nerves which make the iris contract are the ones we have already described, but the nerves which make the iris *expand* come out from no one knows exactly where, go down into the spinal cord back of the chest, into the thoracic sections, out of the spinal cord, up through the neck ganglia, and all the way around and back up into the head in order to run the other end of the iris. In fact, the signal goes through a completely different nervous system, not the central nervous system at all, but the sympathetic nervous system, so it is a very strange way of making things go.

We have already emphasized another strange thing about the eye, that the light-sensitive cells are on the wrong side, so that the light has to go through several layers of other cells before it gets to the receptors—it is built inside out! So some of the features are wonderful and some are apparently stupid.

Figure 36-4 shows the connections of the eye to the part of the brain which is most directly concerned with the visual process. The optic nerve fibers run into a certain area just beyond *D*, called the lateral geniculate, whereupon they run out to a section of the brain called the visual cortex. Notice that some of the fibers from each eye are sent over to the other side of the brain, so the picture formed is incomplete. The optic nerves from the left side of the right eye run across the optic chiasma *B*, while the ones on the left side of the left eye come around and go this same way. So the left side of the brain receives all the information which comes from the left side of the eyeball of each eye, i.e., on the right side of the visual field, while the right side of the brain sees the left side of the visual field. This is the manner in which the information from each of the two eyes is put together in order to tell how far away things are. This is the system of binocular vision.

The connections between the retina and the visual cortex are interesting. If a spot in the retina is excised or destroyed in any way, then the whole fiber will die, and we can thereby find out where it is connected. It turns out that, essentially, the connections are one to one—for each spot in the retina there is one spot in the visual cortex—and spots that are very close together in the retina are very close together in the visual cortex. So the visual cortex still represents the spatial arrangement of the rods and cones, but of course much distorted. Things which are in the center of the field, which occupy a very small part of the retina, are expanded over many, many cells in the visual cortex. It is clear that it is useful to have things which are originally close together, still close together. The most remarkable aspect of the matter, however, is the following. The place where one would think it would be most important to have things close together would be right in the middle of the visual field. Believe it or not, the up-and-down line in our visual field as we look at something is of such a nature that the information from all the points on the right side of that line is going into the left side of the brain, and information from the points on the left side is going into the right side of the brain, and the way this area is made, there is a cut right down through the middle, so that the things that are very close together right in the middle are very far apart in the brain! Somehow, the information has to go from one side of the brain to the other through some other channels, which is quite surprising.

The question of how this network ever gets “wired” together is very interesting. The problem of how much is already wired and how much is learned is an old one. It used to be thought long ago that perhaps it does not have to be wired carefully at all, it is only just roughly interconnected, and then, by experience, the young child learns that when a thing is “up there” it produces some sensation in the brain. (Doctors always tell us what the young child “feels,” but how do *they* know what a child feels at the age of one?) The child, at the age of one, supposedly sees that an object is “up there,” gets a certain sensation, and learns to reach “there,” because when he reaches “here,” it does not work. That approach probably is not correct, because we already see that in many cases there are these special detailed interconnections. More illuminating are some most remarkable experiments done with a salamander. (Incidentally, with the salamander there is a direct crossover connection, without the optic chiasma, because the eyes are on each side of the head and have no common area. Salamanders do not have binocular vision.) The experiment is this. We can cut the optic nerve in a salamander and the nerve will grow out from the eyes again. Thousands and thousands of cell fibers will thus re-establish themselves. Now, in the optic nerve the fibers do not stay adjacent to each other—it is like a great, sloppily made telephone cable, all the fibers twisting and turning, but when it gets to the brain they are all sorted out again. When we cut the optic nerve of the salamander, the interesting question is, will it ever get straightened out? The answer is remarkable: yes. If we cut the optic nerve of the salamander and it grows back, the salamander has good visual acuity again. However, if we cut the optic nerve and *turn the eye upside down* and let it grow back again, it has good visual acuity all right, but it has a terrible error: when the salamander sees a fly “up here,” it jumps at it “down there,” and it never learns. Therefore there is some mysterious way by which the thousands and thousands of fibers find their right places in the brain.

This problem of how much is wired in, and how much is not, is an important problem in the theory of the development of creatures. The answer is not known, but is being studied intensively.

The same experiment in the case of a goldfish shows that there is a terrible knot, like a great scar or complication, in the optic nerve where we cut it, but in spite of all this the fibers grow back to their right places in the brain.

In order to do this, as they grow into the old channels of the optic nerve they must make several decisions about the direction in which they should grow. How do they do this? There seem to be chemical clues that different fibers respond to differently. Think of the enormous number of growing fibers, each of which is an individual differing in some way from its neighbors; in responding to whatever the chemical clues are, it responds in a unique enough way to find its proper place for ultimate connection in the brain! This is an interesting—a fantastic—thing. It is one of the great recently discovered phenomena of biology and is undoubtedly connected to many older unsolved problems of growth, organization, and development of organisms, and particularly of embryos.

One other interesting phenomenon has to do with the motion of the eye. The eyes must be moved in order to make the two images coincide in different circumstances. These motions are of different kinds: one is to follow something, which requires that both eyes must go in the same direction, right or left, and the other is to point them toward the same place at various distances away, which requires that they must move oppositely. The nerves going into the muscles of the eye are already wired up for just such purposes. There is one set of nerves which will pull the muscles on the inside of one eye and the outside of the other, and relax the opposite muscles, so that the two eyes move together. There is another center where an excitation will cause the eyes to move in toward each other from parallel. Either eye can be turned out to the corner if the other eye moves toward the nose, but it is impossible consciously or unconsciously to turn both eyes *out* at the same time, not because there are no *muscles*, but because there is no way to send a signal to turn both eyes out, unless we have had an accident or there is something the matter, for instance if a nerve has been cut. Although the muscles of one eye can certainly steer that eye about, not even a Yogi is able to move *both* eyes out freely

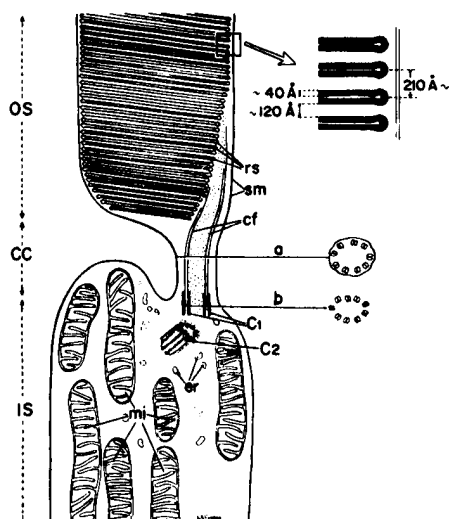


Fig. 36-5. Electron micrograph of a rod cell.

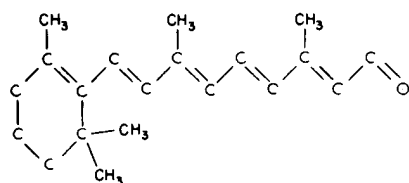


Fig. 36-6. The structure of retinene.

under voluntary control, because there does not seem to be any way to do it. We are already wired to a certain extent. This is an important point, because most of the earlier books on anatomy and psychology, and so on, do not appreciate or do not emphasize the fact that we are so completely wired already—they say that everything is just learned.

36-3 The rod cells

Let us now examine in more detail what happens in the rod cells. Figure 36-5 shows an electron micrograph of the middle of a rod cell (the rod cell keeps going up out of the field). There are layer after layer of plane structures, shown magnified at the right, which contain the substance rhodopsin (visual purple), the dye, or pigment, which produces the effects of vision in the rods. The rhodopsin, which is the pigment, is a big protein which contains a special group called retinene, which can be taken off the protein, and which is, undoubtedly, the main cause of the absorption of light. We do not understand the reason for the planes, but it is very likely that there is some reason for holding all the rhodopsin molecules parallel. The chemistry of the thing has been worked out to a large extent, but there might be some physics to it. It may be that all of the molecules are arranged in some kind of a row so that when one is excited an electron which is generated, say, may run all the way down to some place at the end to get the signal out, or something. This subject is very important, and has not been worked out. It is a field in which both biochemistry and solid state physics, or something like it, will ultimately be used.

This kind of a structure, with layers, appears in other circumstances where light is important, for example in the chloroplast in plants, where the light causes photosynthesis. If we magnify those, we find the same thing with almost the same kind of layers, but there we have chlorophyll, of course, instead of retinene. The chemical form of retinene is shown in Fig. 36-6. It has a series of alternate double bonds along the side chain, which is characteristic of nearly all strongly absorbing organic substances, like chlorophyll, blood, and so on. This substance is impossible for human beings to manufacture in their own cells—we have to eat it. So we eat it in the form of a special substance, which is exactly the same as retinene except that there is a hydrogen tied on the right end; it is called vitamin A, and if we do not eat enough of it, we do not get a supply of retinene, and the eye becomes what we call *night blind*, because there is then not enough pigment in the rhodopsin to see with the rods at night.

The reason why such a series of double bonds absorbs light very strongly is also known. We may just give a hint: The alternating series of double bonds is called a *conjugated* double bond; a double bond means that there is an extra electron there, and this extra electron is easily shifted to the right or left. When light strikes this molecule, the electron of each double bond is shifted over by one step. All the electrons in the whole chain shift, like a string of dominoes falling over, and though each one moves only a little distance (we would expect that, in a single atom, we could move the electron only a little distance), the net effect is the same as though the one at the end was moved over to the other end! It is the same as though one electron went the whole distance back and forth, and so, in this manner, we get a much stronger absorption under the influence of the electric field, than if we could only move the electron a distance which is associated with one atom. So, since it is easy to move the electrons back and forth, retinene absorbs light very strongly; that is the machinery of the physical-chemical end of it.

36-4 The compound (insect) eye

Let us now return to biology. The human eye is not the only kind of eye. In the vertebrates, almost all eyes are essentially like the human eye. However, in the lower animals there are many other kinds of eyes: eye spots, various eye cups, and other less sensitive things, which we have no time to discuss. But there is one other highly developed eye among the invertebrates, the *compound* eye of the insect. (Most insects having large compound eyes also have various additional simpler eyes as well.) A bee is an insect whose vision has been studied very carefully. It is

easy to study the properties of the vision of bees because they are attracted to honey, and we can make experiments in which we identify the honey by putting it on blue paper or red paper, and see which one they come to. By this method some very interesting things have been discovered about the vision of the bee.

In the first place, in trying to measure how acutely bees could see the color difference between two pieces of "white" paper, some researchers found they were not very good, and others found they were fantastically good. Even if the two pieces of white paper were almost exactly the same, the bees could still tell the difference. The experimenters used zinc white for one piece of paper and lead white for the other, and although these look exactly the same to us, the bee could easily distinguish them, because they reflect a different amount in the ultraviolet. In this way it was discovered that the bee's eye is sensitive over a wider range of the spectrum than is our own. Our eye works from 7000 angstroms to 4000 angstroms, from red to violet, but the bee's can see down to 3000 angstroms into the ultraviolet! This makes for a number of different interesting effects. In the first place, bees can distinguish between many flowers which to us look alike. Of course, we must realize that the colors of flowers are not designed for *our* eyes, but for the bee; they are signals to attract the bees to a specific flower. We all know that there are many "white" flowers. Apparently white is not very interesting to the bees, because it turns out that all of the white flowers have different proportions of reflection in the *ultraviolet*; they do not reflect one hundred percent of the ultraviolet as would a true white. All the light is not coming back, the ultraviolet is missing, and that is a color, just as, for us, if the blue is missing, it comes out yellow. So, all the flowers are colored for the bees. However, we also know that red cannot be seen by bees. Thus we might expect that all red flowers should look black to the bee. Not so! A careful study of red flowers shows, first, that even with our own eye we can see that a great majority of red flowers have a bluish tinge because they are mainly reflecting an additional amount in the blue, which is the part that the bee sees. Furthermore, experiments also show that flowers vary in their reflection of the ultraviolet over different parts of the petals, and so on. So if we could see the flowers as bees see them they would be even more beautiful and varied!

It has been shown, however, that there are a few red flowers which do *not* reflect in the blue or in the ultraviolet, and *would*, therefore, appear black to the bee! This was of quite some concern to the people who worry about this matter, because black does not seem like an interesting color, since it is hard to tell from a dirty old shadow. It actually turned out that these flowers were *not* visited by bees, these are the flowers that are visited by *hummingbirds*, and hummingbirds *can* see the red!

Another interesting aspect of the vision of the bee is that bees can apparently tell the direction of the sun by looking at a patch of blue sky, without seeing the sun itself. We cannot easily do this. If we look out the window at the sky and see that it is blue, in which direction is the sun? The bee can tell, because the bee is quite sensitive to the *polarization* of light, and the scattered light of the sky is polarized.* There is still some debate about how this sensitivity operates. Whether it is because the reflections of the light are different in different circumstances, or the bee's eye is directly sensitive, is not yet known.†

It is also said that the bee can notice flicker up to 200 oscillations per second, while we see it only up to 20. The motions of bees in the hives are very quick; the feet move and the wings vibrate, but it is very hard for us to see these motions with our eye. However, if we could see more rapidly we would be able to see the motion. It is probably very important to the bee that its eye has such a rapid response.

* The human eye also has a slight sensitivity to the polarization of light, and one *can* learn to tell the direction of the sun! The phenomenon that is involved here is called *Haidinger's brush*; it is a faint, yellowish hourglass-like pattern seen at the center of the visual field when one looks at a broad, featureless expanse using polarizing glasses. It can also be seen in the blue sky without polarizing glasses if one rotates his head back and forth about the axis of vision.

† Evidence obtained since this lecture was given indicates that the eye is directly sensitive.



Fig. 36-7. The structure of an ommatidium (a single cell of a compound eye).

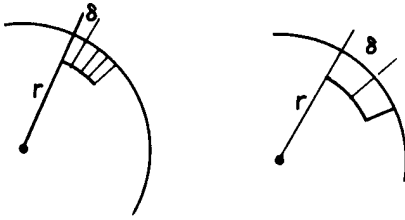


Fig. 36-8. Schematic view of packing of ommatidia in the eye of a bee.

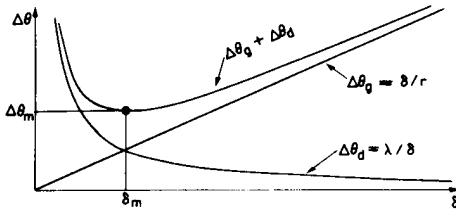


Fig. 36-9. The optimum size for an ommatidium is δ_m .

Now let us discuss the visual acuity we could expect from the bee. The eye of a bee is a compound eye, and it is made of a large number of special cells called *ommatidia*, which are arranged conically on the surface of a sphere (roughly) on the outside of the bee's head. Figure 36-7 shows a picture of one such ommatidium. At the top there is a transparent area, a kind of "lens," but actually it is more like a filter or light pipe to make the light come down along the narrow fiber, which is where the absorption presumably occurs. Out of the other end of it comes the nerve fiber. The central fiber is surrounded on its sides by six cells which, in fact, have secreted the fiber. That is enough description for our purposes; the point is that it is a conical thing and many can fit next to each other all over the surface of the eye of the bee.

Now let us discuss the resolution of the eye of the bee. If we draw lines (Fig. 36-8) to represent the ommatidia on the surface, which we suppose is a sphere of radius r , we may actually *calculate* how wide each ommatidium is by using our brains, and assuming that evolution is as clever as we are! If we have a very large ommatidium we do not have much resolution. That is, one cell gets a piece of information from one direction, and the adjacent cell gets a piece of information from another direction, and so on, and the bee cannot see things in between very well. So the uncertainty of visual acuity in the eye will surely correspond to an angle, the angle of the end of the ommatidium relative to the center of curvature of the eye. (The eye cells, of course, exist only at the surface of the sphere; inside that is the head of the bee.) This angle, from one ommatidium to the next, is, of course, the diameter of the ommatidia divided by the radius of the eye surface:

$$\Delta\theta_g = \delta/r. \quad (36.1)$$

So, we may say, "The finer we make the δ , the more the visual acuity. So why doesn't the bee just use very, very fine ommatidia?" *Answer:* We know enough physics to realize that if we are trying to get light down into a narrow slot, we cannot see accurately in a given direction because of the diffraction effect. The light that comes from several directions can enter and, due to diffraction, we will get light coming in at angle $\Delta\theta_d$ such that

$$\Delta\theta_d = \lambda/\delta. \quad (36.2)$$

Now we see that if we make the δ too small, then each ommatidium does not look in only one direction, because of diffraction! If we make them too big, each one sees in a definite direction, but there are not enough of them to get a good view of the scene. So we adjust the distance d in order to make minimal the total effect of these two. If we add the two together, and find the place where the sum has a minimum (Fig. 36-9), we find that

$$\frac{d(\Delta\theta_g + \Delta\theta_d)}{d\delta} = 0 = \frac{1}{r} - \frac{\lambda}{\delta^2}, \quad (36.3)$$

which gives us a distance

$$\delta = \sqrt{\lambda r}. \quad (36.4)$$

If we guess that r is about 3 millimeters, take the light that the bee sees as 4000 angstroms, and put the two together and take the square root, we find

$$\begin{aligned} \delta &= (3 \times 10^{-3} \times 4 \times 10^{-7})^{1/2} \text{ m} \\ &= 3.5 \times 10^{-5} \text{ m} = 35 \mu. \end{aligned} \quad (36.5)$$

The book says the diameter is 30μ , so that is rather good agreement! So, apparently, it really works, and we can understand what determines the size of the bee's eye! It is also easy to put the above number back in and find out how good the bee's eye actually is in angular resolution; it is very poor relative to our own. We can see things that are thirty times smaller in apparent size than the bee; the bee has a rather fuzzy out-of-focus image relative to what we can see. Nevertheless it is all right, and it is the best they can do. We might ask why the bees do not develop

a good eye like our own, with a lens and so on. There are several interesting reasons. In the first place, the bee is too small; if it had an eye like ours, but on his scale, the opening would be about $30\ \mu$ in size and diffraction would be so important that it would not be able to see very well anyway. The eye is not good if it is too small. Secondly, if it were as big as the bee's head, then the eye would occupy the whole head of the bee. The beauty of the compound eye is that it takes up no space, it is just a very thin layer on the surface of the bee. So when we argue that they should have done it our way, we must remember that they had their own problems!

36-5 Other eyes

Besides the bees, many other animals can see color. Fish, butterflies, birds, and reptiles can see color, but it is believed that most mammals cannot. The primates can see color. The birds certainly see color, and that accounts for the colors of birds. There would be no point in having such brilliantly colored males if the females could not notice it! That is, the evolution of the sexual "whatever it is" that the birds have is a result of the female being able to see color. So next time we look at a peacock and think of what a brilliant display of gorgeous color it is, and how delicate all the colors are, and what a wonderful aesthetic sense it takes to appreciate all that, we should not compliment the peacock, but should compliment the visual acuity and aesthetic sense of the *peahen*, because that is what has generated the beautiful scene!

All invertebrates have poorly developed eyes or compound eyes, but all the vertebrates have eyes very similar to our own, with one exception. If we consider the highest form of animal, we usually say, "Here we are!," but if we take a less prejudiced point of view and restrict ourselves to the invertebrates, so that we cannot include ourselves, and ask what is the highest invertebrate animal, most zoologists agree that the *octopus* is the highest animal! It is very interesting that, besides the development of its brain and its reactions and so on, which are rather good for an invertebrate, it has also developed, independently, a different eye. It is not a compound eye or an eye spot—it has a cornea, it has lids, it has an iris, it has a lens, it has two regions of water, it has a retina behind. It is essentially the same as the eye of the vertebrates! It is a remarkable example of a coincidence in evolution where nature has twice discovered the same solution to a problem, with one slight improvement. In the octopus it also turns out, amazingly, that the retina is a piece of the brain that has come out in the same way in its embryonic development as is true for vertebrates, but the interesting thing which is different is that the cells which are sensitive to light are on the *inside*, and the cells which do the calculation are in back of them, rather than "inside out," as in our eye. So we see, at least, that there is no good reason for its being inside out. The other time nature tried it, she got it straightened out! (See Fig. 36-10.) The biggest eyes in the world are those of the giant squid; they have been found up to 15 inches in diameter!

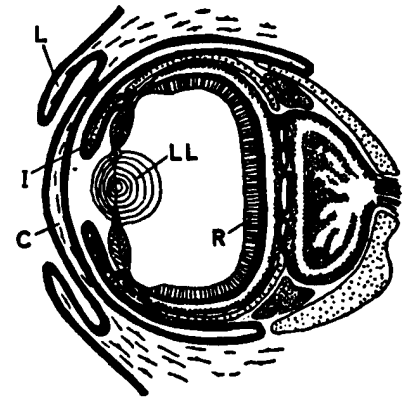


Fig. 36-10. The eye of an octopus.

36-6 Neurology of vision

One of the main points of our subject is the interconnection of information from one part of the eye to the other. Let us consider the compound eye of the horseshoe crab, on which considerable experimentation has been done. First of all, we must appreciate what kind of information can come along nerves. A nerve carries a kind of disturbance which has an electrical effect that is easy to detect, a kind of wavelike disturbance which runs down the nerve and produces an effect at the other end: a long piece of the nerve cell, called the axon, carries the information along, and a certain kind of impulse, called a "spike," goes along if it is excited at one end. When one spike goes down the nerve, another cannot immediately follow. All the spikes are of the same size, so it is not that we get *higher* spikes when the thing is more strongly excited, but that we get *more spikes per second*. The *size* of the spike is determined by the fiber. It is important to appreciate this in order to see what happens next.

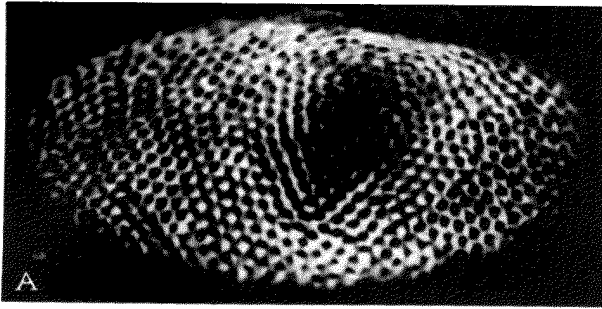


Fig. 36-11. The compound eye of the horseshoe crab. (a) Normal view. (b) Cross section.

Figures 36-7, 11, 12, 13 reprinted with permission from Goldsmith, *Sensory Communications*, W. A. Rosenblith, ed. Copyright 1961, Massachusetts Institute of Technology.

Figure 36-11(a) shows the compound eye of the horseshoe crab; it is not very much of an eye, it has only about a thousand ommatidia. Figure 36-11(b) is a cross section through the system; one can see the ommatidia, with the nerve fibers that run out of them and go into the brain. But note that even in a horseshoe crab there are little interconnections. They are much less elaborate than in the human eye, and it gives us a chance to study a simpler example.

Let us now look at the experiments which have been done by putting fine electrodes into the optic nerve of the horseshoe crab, and shining light on only one of the ommatidia, which is easy to do with lenses. If we turn a light on at some instant t_0 , and measure the electrical pulses that come out, we find that there is a slight delay and then a rapid series of discharges which gradually slow down to a uniform rate, as shown in Fig. 36-12(a). When the light goes out, the discharge stops. Now it is very interesting that if, while our amplifier is connected to this same nerve fiber, we shine light on a *different* ommatidium nothing happens; no signal.

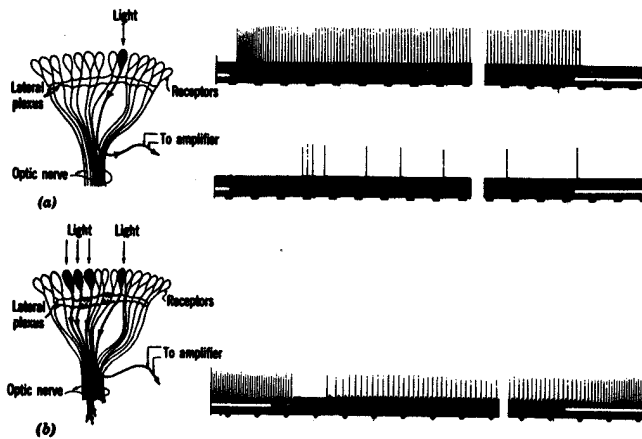


Fig. 36-12. The response to light of the nerve fibers of the eye of the horseshoe crab.

Now we do another experiment: we shine the light on the original ommatidium and get the same response, but if we now turn light on another one nearby as well, the pulses are interrupted briefly and then run at a much lower rate (Fig. 36-12b). The rate of one is inhibited by the impulses which are coming out of the other! In other words, each nerve fiber carries the information from one ommatidium, but the amount that it carries is inhibited by the signals from the others. So, for example, if the whole eye is more or less uniformly illuminated, the information coming from any one ommatidium will be relatively weak, because it is inhibited by so many. In fact the inhibition is additive—if we shine light on several nearby ommatidia the inhibition is very great. The inhibition is greater when the ommatidia are closer, and if the ommatidia are far enough away from one another, inhibition is practically zero. So it is additive and depends on the distance; here is a first example of information from different parts of the eye being combined in the eye itself. We can see, perhaps, if we think about it awhile, that this is a device to *enhance contrast* at the edges of objects, because if a part of the scene is light and a part is black, then the ommatidia in the lighted area give impulses that are

inhibited by all the other light in the neighborhood, so it is relatively weak. On the other hand, an ommatidium at the boundary which is given a “white” impulse is also inhibited by others in the neighborhood, but there are not as many of them, since some are black; the net signal is therefore stronger. The result would be a curve, something like that of Fig. 36-13. The crab will see an enhancement of the contour.

The fact that there is an enhancement of contours has long been known; in fact it is a remarkable thing that has been commented on by psychologists many times. In order to draw an object, we have only to draw its outline. How used we are to looking at pictures that have only the outline! What is the outline? The outline is only the edge difference between light and dark or one color and another. It is not something definite. It is not, believe it or not, that every object has a line around it! There is no such line. It is only in our own psychological makeup that there is a line; we are beginning to understand the reasons why the “line” is enough of a clue to get the whole thing. Presumably our own eye works in some similar manner—much more complicated, but similar.

Finally, we shall briefly describe the more elaborate work, the beautiful, advanced work that has been done on the frog. Doing a corresponding experiment on a frog, by putting very fine, beautifully built needlelike probes into the optic nerve of a frog, one can obtain the signals that are going along one particular axon and, just as in the case of the horseshoe crab, we find that the information does not depend on just one spot in the eye, but is a sum of information over several spots.

The most recent picture of the operation of the frog’s eye is the following. One can find four different kinds of optic nerve fibers, in the sense that there are four different kinds of responses. These experiments were not done by shining on-and-off impulses of light, because that is not what a frog sees. A frog just sits there and his eyes never move, unless the lily pad is flopping back and forth, and in that case his eyes wobble just right so that the image stays put. He does not turn his eyes. If anything moves in his field of vision, like a little bug (he has to be able to see something small moving in the fixed background), it turns out that there are four different kinds of fibers which discharge, whose properties are summarized in Table 36-1. Sustained edge detection, nonerasable, means that if we bring an object with an edge into the field of view of the frog, then there are a lot of impulses in this particular fiber while the object is moving, but they die down to a sustained impulse that continues as long as the edge is there, even if it is standing still. If we turn out the light, the impulses stop. If we turn it on again while the edge is still in view, they start again. They are not erasable. Another kind of fiber is very similar, except that if the edge is straight it does not work. It must be a convex edge with dark behind it! How complicated must be the system of interconnections in the retina of the eye of the frog in order for it to understand that a convex surface has moved in! Furthermore, although this fiber does sustain somewhat, it does not sustain as long as the other, and if we turn out the light and turn it on again it does *not* build up again. It depends on the moving in of the convex

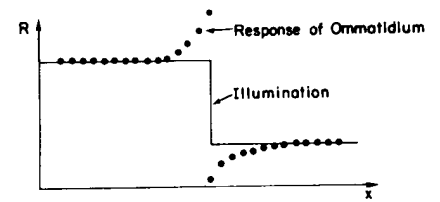


Fig. 36-13. The net response of horseshoe crab ommatidia near a sharp change in illumination.

Table 36-1

Types of response in optic nerve fibers of a frog

Type	Speed	Angular field
1. Sustained edge detection (nonerasable)	0.2–0.5 m/sec	1°
2. Convex edge detection (erasable)	0.5 m/sec	2°–3°
3. Changing contrast detection	1–2 m/sec	7°–10°
4. Dimming detection	Up to $\frac{1}{2}$ m/sec	Up to 15°
5. Darkness detection	?	Very large

surface. The eye sees it move in and remembers that it is there, but if we merely turn out the light for a moment, it simply forgets it and no longer sees it.

Another example is change-in-contrast detection. If there is an edge moving in or out there are pulses, but if the thing stands still there are no pulses at all.

Then there is a dimming detector. If the light intensity is going down it creates pulses, but if it stays down or stays up, the impulse stops; it only works while the light is dimming.

Then, finally, there are a few fibers which are dark detectors—a most amazing thing—they fire all the time! If we increase the light, they fire less rapidly, but all the time. If we decrease the light, they fire more rapidly, all the time. In the dark they fire like mad, perpetually saying, “It is dark! It is dark! It is dark!”

Now these responses seem to be rather complicated to classify, and we might wonder whether perhaps the experiments are being misinterpreted. But it is very interesting that these same classes are very clearly separated in the anatomy of the frog! By other measurements, after these responses had been classified (*afterwards*, that is what is important about this), it was discovered that the *speed* of the signals on the different fibers was not the same, so here was another, independent way to check which kind of a fiber we have found!

Another interesting question is from how big an area is one particular fiber making its calculations? The answer is different for the different classes.

Figure 36-14 shows the surface of the so-called tectum of a frog, where the nerves come into the brain from the optic nerve. All the nerve fibers coming in from the optic nerve make connections in various layers of the tectum. This layered structure is analogous to the retina; that is partly why we know that the brain and retina are very similar. Now, by taking an electrode and moving it down in succession through the layers, we can find out which kinds of optic nerves end where, and the beautiful and wonderful result is that the different kinds of fibers end in different layers! The first ones end in number 1 type, the second in number 2, the threes and fives end in the same place, and deepest of all is number four. (What a coincidence, they got the numbers almost in the right order! No, that is why they numbered them that way, the first paper had the numbers in a different order!)

We may briefly summarize what we have just learned this way: There are three pigments, presumably. There may be many different kinds of receptor cells containing the three pigments in different proportions, but there are many cross connections which may permit additions and subtractions through addition and reinforcement in the nervous system. So before we really understand color vision, we will have to understand the final sensation. This subject is still an open one, but these researches with microelectrodes and so on will perhaps ultimately give us more information on how we see color.



Fig. 36-14. The tectum of a frog.

BIBLIOGRAPHY

Committee on Colorimetry, Optical Society of America, *The Science of Color*, Thomas Y. Crowell Company, New York, 1953.

“Mechanisms of Vision,” 2nd Supplement to *Journal of General Physiology*, Vol. 43, No. 6, Part 2, July 1960, Rockefeller Institute Press.

SPECIFIC ARTICLES:

DEROBERTIS, E., “Some Observations on the Ultrastructure and Morphogenesis of Photoreceptors,” pp. 1-15.

HURVICH, L. M. and D. JAMESON, “Perceived Color, Induction Effects, and Opponent-Response Mechanisms,” pp. 63-80.

ROSENBLITH, W. A., ed., *Sensory Communication*, Massachusetts Institute of Technology Press, Cambridge, Mass., 1961.

“Sight, Sense of,” *Encyclopaedia Britannica*, Vol. 20, 1957, pp. 628-635.

Quantum Behavior

37-1 Atomic mechanics

In the last few chapters we have treated the essential ideas necessary for an understanding of most of the important phenomena of light—or electromagnetic radiation in general. (We have left a few special topics for next year. Specifically, the theory of the index of dense materials and total internal reflection.) What we have dealt with is called the “classical theory” of electric waves, which turns out to be a completely adequate description of nature for a large number of effects. We have not had to worry yet about the fact that light energy comes in lumps or “photons.”

We would like to take up as our next subject the problem of the behavior of relatively large pieces of matter—their mechanical and thermal properties, for instance. In discussing these, we will find that the “classical” (or older) theory fails almost immediately, because matter is really made up of atomic-sized particles. Still, we will deal only with the classical part, because that is the only part that we can understand using the classical mechanics we have been learning. But we shall not be very successful. We shall find that in the case of matter, unlike the case of light, we shall be in difficulty relatively soon. We could, of course, continuously skirt away from the atomic effects, but we shall instead interpose here a short excursion in which we will describe the basic ideas of the quantum properties of matter, i.e., the quantum ideas of atomic physics, so that you will have some feeling for what it is we are leaving out. For we will have to leave out some important subjects that we cannot avoid coming close to.

So we will give now the *introduction* to the subject of quantum mechanics, but will not be able actually to get into the subject until much later.

“Quantum mechanics” is the description of the behavior of matter in all its details and, in particular, of the happenings on an atomic scale. Things on a very small scale behave like nothing that you have any direct experience about. They do not behave like waves, they do not behave like particles, they do not behave like clouds, or billiard balls, or weights on springs, or like anything that you have ever seen.

Newton thought that light was made up of particles, but then it was discovered, as we have seen here, that it behaves like a wave. Later, however (in the beginning of the twentieth century) it was found that light did indeed sometimes behave like a particle. Historically, the electron, for example, was thought to behave like a particle, and then it was found that in many respects it behaved like a wave. So it really behaves like neither. Now we have given up. We say: “It is like *neither*.”

There is one lucky break, however—electrons behave just like light. The quantum behavior of atomic objects (electrons, protons, neutrons, photons, and so on) is the same for all, they are all “particle waves,” or whatever you want to call them. So what we learn about the properties of electrons (which we shall use for our examples) will apply also to all “particles,” including photons of light.

The gradual accumulation of information about atomic and small-scale behavior during the first quarter of this century, which gave some indications about how small things do behave, produced an increasing confusion which was finally resolved in 1926 and 1927 by Schrödinger, Heisenberg, and Born. They finally obtained a consistent description of the behavior of matter on a small scale. We take up the main features of that description in this chapter.

Because atomic behavior is so unlike ordinary experience, it is very difficult to get used to and it appears peculiar and mysterious to everyone, both to the

37-1 Atomic mechanics

37-2 An experiment with bullets

37-3 An experiment with waves

37-4 An experiment with electrons

37-5 The interference of electron waves

37-6 Watching the electrons

37-7 First principles of quantum mechanics

37-8 The uncertainty principle

novice and to the experienced physicist. Even the experts do not understand it the way they would like to, and it is perfectly reasonable that they should not, because all of direct, human experience and of human intuition applies to large objects. We know how large objects will act, but things on a small scale just do not act that way. So we have to learn about them in a sort of abstract or imaginative fashion and not by connection with our direct experience.

In this chapter we shall tackle immediately the basic element of the mysterious behavior in its most strange form. We choose to examine a phenomenon which is impossible, *absolutely* impossible, to explain in any classical way, and which has in it the heart of quantum mechanics. In reality, it contains the *only* mystery. We cannot explain the mystery in the sense of “explaining” how it works. We will *tell* you how it works. In telling you how it works we will have told you about the basic peculiarities of all quantum mechanics.

37-2 An experiment with bullets

To try to understand the quantum behavior of electrons, we shall compare and contrast their behavior, in a particular experimental setup, with the more familiar behavior of particles like bullets, and with the behavior of waves like water waves. We consider first the behavior of bullets in the experimental setup shown diagrammatically in Fig. 37-1. We have a machine gun that shoots a stream of bullets. It is not a very good gun, in that it sprays the bullets (randomly) over a fairly large angular spread, as indicated in the figure. In front of the gun we have a wall (made of armor plate) that has in it two holes just about big enough to let a bullet through. Beyond the wall is a backstop (say a thick wall of wood) which will “absorb” the bullets when they hit it. In front of the wall we have an object which we shall call a “detector” of bullets. It might be a box containing sand. Any bullet that enters the detector will be stopped and accumulated. When we wish, we can empty the box and count the number of bullets that have been caught. The detector can be moved back and forth (in what we will call the x -direction). With this apparatus, we can find out experimentally the answer to the question: “What is the probability that a bullet which passes through the holes in the wall will arrive at the backstop at the distance x from the center?” First, you should realize that we should talk about probability, because we cannot say definitely where any particular bullet will go. A bullet which happens to hit one of the holes may bounce off the edges of the hole, and may end up anywhere at all. By “probability” we mean the chance that the bullet will arrive at the detector, which we can measure by counting the number which arrive at the detector in a certain time and then taking the ratio of this number to the *total* number that hit the backstop during that time. Or, if we assume that the gun always shoots at the same rate during the measurements, the probability we want is just proportional to the number that reach the detector in some standard time interval.

For our present purposes we would like to imagine a somewhat idealized experiment in which the bullets are not real bullets, but are *indestructible* bullets—they cannot break in half. In our experiment we find that bullets always arrive in lumps, and when we find something in the detector, it is always one whole bullet. If the rate at which the machine gun fires is made very low, we find that at any given moment either nothing arrives, or one and only one—exactly one—bullet arrives at the backstop. Also, the size of the lump certainly does not depend on the rate of firing of the gun. We shall say: “Bullets *always* arrive in identical lumps.” What we measure with our detector is the probability of arrival of a lump. And we measure the probability as a function of x . The result of such measurements with this apparatus (we have not yet done the experiment, so we are really imagining the result) are plotted in the graph drawn in part (c) of Fig. 37-1. In the graph we plot the probability to the right and x vertically, so that the x -scale fits the diagram of the apparatus. We call the probability P_{12} because the bullets may have come either through hole 1 or through hole 2. You will not be surprised that P_{12} is large near the middle of the graph but gets small if x is very large. You may wonder, however, why P_{12} has its maximum value at $x = 0$. We can understand

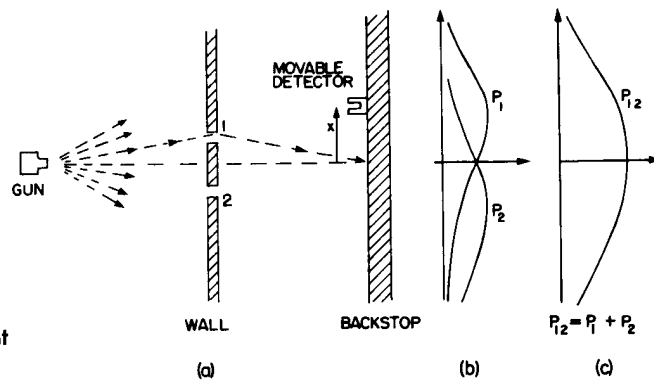


Fig. 37-1. Interference experiment with bullets.

this fact if we do our experiment again after covering up hole 2, and once more while covering up hole 1. When hole 2 is covered, bullets can pass only through hole 1, and we get the curve marked P_1 in part (b) of the figure. As you would expect, the maximum of P_1 occurs at the value of x which is on a straight line with the gun and hole 1. When hole 1 is closed, we get the symmetric curve P_2 drawn in the figure. P_2 is the probability distribution for bullets that pass through hole 2. Comparing parts (b) and (c) of Fig. 37-1, we find the important result that

$$P_{12} = P_1 + P_2. \quad (37.1)$$

The probabilities just add together. The effect with both holes open is the sum of the effects with each hole open alone. We shall call this result an observation of “no interference,” for a reason that you will see later. So much for bullets. They come in lumps, and their probability of arrival shows no interference.

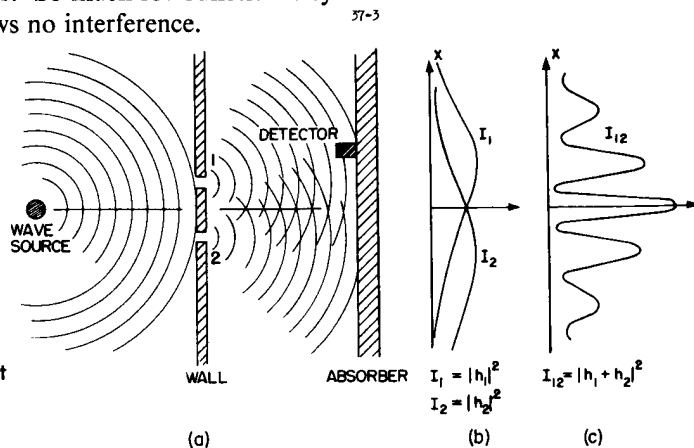


Fig. 37-2. Interference experiment with water waves.

37-3 An experiment with waves

Now we wish to consider an experiment with water waves. The apparatus is shown diagrammatically in Fig. 37-2. We have a shallow trough of water. A small object labeled the “wave source” is jiggled up and down by a motor and makes circular waves. To the right of the source we have again a wall with two holes, and beyond that is a second wall, which, to keep things simple, is an “absorber,” so that there is no reflection of the waves that arrive there. This can be done by building a gradual sand “beach.” In front of the beach we place a detector which can be moved back and forth in the x -direction, as before. The detector is now a device which measures the “intensity” of the wave motion. You can imagine a gadget which measures the height of the wave motion, but whose scale is calibrated in proportion to the *square* of the actual height, so that the reading is proportional to the intensity of the wave. Our detector reads, then, in proportion to the *energy* being carried by the wave—or rather, the rate at which energy is carried to the detector.

With our wave apparatus, the first thing to notice is that the intensity can have *any* size. If the source just moves a very small amount, then there is just a little bit of wave motion at the detector. When there is more motion at the source,

there is more intensity at the detector. The intensity of the wave can have any value at all. We would *not* say that there was any “lumpiness” in the wave intensity.

Now let us measure the wave intensity for various values of x (keeping the wave source operating always in the same way). We get the interesting-looking curve marked I_{12} in part (c) of the figure.

We have already worked out how such patterns can come about when we studied the interference of electric waves. In this case we would observe that the original wave is diffracted at the holes, and new circular waves spread out from each hole. If we cover one hole at a time and measure the intensity distribution at the absorber we find the rather simple intensity curves shown in part (b) of the figure. I_1 is the intensity of the wave from hole 1 (which we find by measuring when hole 2 is blocked off) and I_2 is the intensity of the wave from hole 2 (seen when hole 1 is blocked).

The intensity I_{12} observed when both holes are open is certainly *not* the sum of I_1 and I_2 . We say that there is “interference” of the two waves. At some places (where the curve I_{12} has its maxima) the waves are “in phase” and the wave peaks add together to give a large amplitude and, therefore, a large intensity. We say that the two waves are “interfering constructively” at such places. There will be such constructive interference wherever the distance from the detector to one hole is a whole number of wavelengths larger (or shorter) than the distance from the detector to the other hole.

At those places where the two waves arrive at the detector with a phase difference of π (where they are “out of phase”) the resulting wave motion at the detector will be the difference of the two amplitudes. The waves “interfere destructively,” and we get a low value for the wave intensity. We expect such low values wherever the distance between hole 1 and the detector is different from the distance between hole 2 and the detector by an odd number of half-wavelengths. The low values of I_{12} in Fig. 37-2 correspond to the places where the two waves interfere destructively.

You will remember that the quantitative relationship between I_1 , I_2 , and I_{12} can be expressed in the following way: The instantaneous height of the water wave at the detector for the wave from hole 1 can be written as (the real part of) $\hat{h}_1 e^{i\omega t}$, where the “amplitude” \hat{h}_1 is, in general, a complex number. The intensity is proportional to the mean squared height or, when we use the complex numbers, to $|\hat{h}_1|^2$. Similarly, for hole 2 the height is $\hat{h}_2 e^{i\omega t}$ and the intensity is proportional to $|\hat{h}_2|^2$. When both holes are open, the wave heights add to give the height $(\hat{h}_1 + \hat{h}_2)e^{i\omega t}$ and the intensity $|\hat{h}_1 + \hat{h}_2|^2$. Omitting the constant of proportionality for our present purposes, the proper relations for *interfering waves* are

$$I_1 = |\hat{h}_1|^2, \quad I_2 = |\hat{h}_2|^2, \quad I_{12} = |\hat{h}_1 + \hat{h}_2|^2. \quad (37.2)$$

You will notice that the result is quite different from that obtained with bullets (Eq. 37.1). If we expand $|\hat{h}_1 + \hat{h}_2|^2$ we see that

$$|\hat{h}_1 + \hat{h}_2|^2 = |\hat{h}_1|^2 + |\hat{h}_2|^2 + 2|\hat{h}_1||\hat{h}_2| \cos \delta, \quad (37.3)$$

where δ is the phase difference between \hat{h}_1 and \hat{h}_2 . In terms of the intensities, we could write

$$I_{12} = I_1 + I_2 + 2\sqrt{I_1 I_2} \cos \delta. \quad (37.4)$$

The last term in (37.4) is the “interference term.” So much for water waves. The intensity can have any value, and it shows interference.

37-4 An experiment with electrons

Now we imagine a similar experiment with electrons. It is shown diagrammatically in Fig. 37-3. We make an electron gun which consists of a tungsten wire heated by an electric current and surrounded by a metal box with a hole in it. If the wire is at a negative voltage with respect to the box, electrons emitted by the wire will be accelerated toward the walls and some will pass through the hole. All the electrons which come out of the gun will have (nearly) the same energy. In front of the gun is again a wall (just a thin metal plate) with two holes in it.

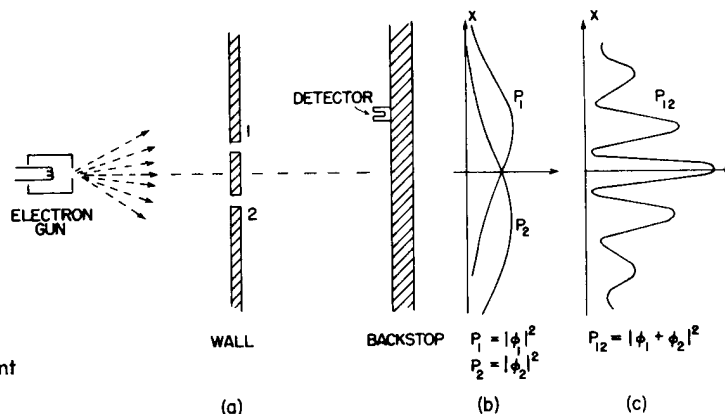


Fig. 37-3. Interference experiment with electrons.

Beyond the wall is another plate which will serve as a “backstop.” In front of the backstop we place a movable detector. The detector might be a geiger counter or, perhaps better, an electron multiplier, which is connected to a loudspeaker.

We should say right away that you should not try to set up this experiment (as you could have done with the two we have already described). This experiment has never been done in just this way. The trouble is that the apparatus would have to be made on an impossibly small scale to show the effects we are interested in. We are doing a “thought experiment,” which we have chosen because it is easy to think about. We know the results that *would* be obtained because there *are* many experiments that have been done, in which the scale and the proportions have been chosen to show the effects we shall describe.

The first thing we notice with our electron experiment is that we hear sharp “clicks” from the detector (that is, from the loudspeaker). And all “clicks” are the same. There are *no* “half-clicks.”

We would also notice that the “clicks” come very erratically. Something like: click click-click . . . click click click-click click . . . , etc., just as you have, no doubt, heard a geiger counter operating. If we count the clicks which arrive in a sufficiently long time—say for many minutes—and then count again for another equal period, we find that the two numbers are very nearly the same. So we can speak of the *average rate* at which the clicks are heard (so-and-so-many clicks per minute on the average).

As we move the detector around, the *rate* at which the clicks appear is faster or slower, but the size (loudness) of each click is always the same. If we lower the temperature of the wire in the gun the rate of clicking slows down, but still each click sounds the same. We would notice also that if we put two separate detectors at the backstop, one *or* the other would click, but never both at once. (Except that once in a while, if there were two clicks very close together in time, our ear might not sense the separation.) We conclude, therefore, that whatever arrives at the backstop arrives in “lumps.” All the “lumps” are the same size: only whole “lumps” arrive, and they arrive one at a time at the backstop. We shall say: “Electrons always arrive in identical lumps.”

Just as for our experiment with bullets, we can now proceed to find experimentally the answer to the question: “What is the relative probability that an electron ‘lump’ will arrive at the backstop at various distances x from the center?” As before, we obtain the relative probability by observing the rate of clicks, holding the operation of the gun constant. The probability that lumps will arrive at a particular x is proportional to the average rate of clicks at that x .

The result of our experiment is the interesting curve marked P_{12} in part (c) of Fig. 37-3. Yes! That is the way electrons go.

37-5 The interference of electron waves

Now let us try to analyze the curve of Fig. 37-3 to see whether we can understand the behavior of the electrons. The first thing we would say is that since they come in lumps, each lump, which we may as well call an electron, has come either through hole 1 or through hole 2. Let us write this in the form of a “Proposition”:

Proposition A: Each electron *either* goes through hole 1 *or* it goes through hole 2.

Assuming Proposition A, all electrons that arrive at the backstop can be divided into two classes: (1) those that come through hole 1, and (2) those that come through hole 2. So our observed curve must be the sum of the effects of the electrons which come through hole 1 and the electrons which come through hole 2. Let us check this idea by experiment. First, we will make a measurement for those electrons that come through hole 1. We block off hole 2 and make our counts of the clicks from the detector. From the clicking rate, we get P_1 . The result of the measurement is shown by the curve marked P_1 in part (b) of Fig. 37-3. The result seems quite reasonable. In a similar way, we measure P_2 , the probability distribution for the electrons that come through hole 2. The result of this measurement is also drawn in the figure.

The result P_{12} obtained with *both* holes open is clearly not the sum of P_1 and P_2 , the probabilities for each hole alone. In analogy with our water-wave experiment, we say: "There is interference."

$$\text{For electrons: } P_{12} \neq P_1 + P_2. \quad (37.5)$$

How can such an interference come about? Perhaps we should say: "Well, that means, presumably, that it is *not true* that the lumps go either through hole 1 or hole 2, because if they did, the probabilities should add. Perhaps they go in a more complicated way. They split in half and . . ." But no! They cannot, they always arrive in lumps . . . "Well, perhaps some of them go through 1, and then they go around through 2, and then around a few more times, or by some other complicated path . . . then by closing hole 2, we changed the chance that an electron that *started out* through hole 1 would finally get to the backstop . . ." But notice! There are some points at which very few electrons arrive when *both* holes are open, but which receive many electrons if we close one hole, so *closing* one hole *increased* the number from the other. Notice, however, that at the center of the pattern, P_{12} is more than twice as large as $P_1 + P_2$. It is as though closing one hole *decreased* the number of electrons which come through the other hole. It seems hard to explain *both* effects by proposing that the electrons travel in complicated paths.

It is all quite mysterious. And the more you look at it the more mysterious it seems. Many ideas have been concocted to try to explain the curve for P_{12} in terms of individual electrons going around in complicated ways through the holes. None of them has succeeded. None of them can get the right curve for P_{12} in terms of P_1 and P_2 .

Yet, surprisingly enough, the *mathematics* for relating P_1 and P_2 to P_{12} is extremely simple. For P_{12} is just like the curve I_{12} of Fig. 37-2, and *that* was simple. What is going on at the backstop can be described by two complex numbers that we can call $\hat{\phi}_1$ and $\hat{\phi}_2$ (they are functions of x , of course). The absolute square of $\hat{\phi}_1$ gives the effect with only hole 1 open. That is, $P_1 = |\hat{\phi}_1|^2$. The effect with only hole 2 open is given by $\hat{\phi}_2$ in the same way. That is, $P_2 = |\hat{\phi}_2|^2$. And the combined effect of the two holes is just $P_{12} = |\hat{\phi}_1 + \hat{\phi}_2|^2$. The *mathematics* is the same as that we had for the water waves! (It is hard to see how one could get such a simple result from a complicated game of electrons going back and forth through the plate on some strange trajectory.)

We conclude the following: The electrons arrive in lumps, like particles, and the probability of arrival of these lumps is distributed like the distribution of intensity of a wave. It is in this sense that an electron behaves "sometimes like a particle and sometimes like a wave."

Incidentally, when we were dealing with classical waves we defined the intensity as the mean over time of the square of the wave amplitude, and we used complex numbers as a mathematical trick to simplify the analysis. But in quantum mechanics it turns out that the amplitudes *must* be represented by complex numbers. The real parts alone will not do. That is a technical point, for the moment, because the formulas look just the same.

Since the probability of arrival through both holes is given so simply, although it is not equal to $(P_1 + P_2)$, that is really all there is to say. But there are a large number of subtleties involved in the fact that nature does work this way. We would like to illustrate some of these subtleties for you now. First, since the number that arrives at a particular point is *not* equal to the number that arrives through 1 plus the number that arrives through 2, as we would have concluded from Proposition A, undoubtedly we should conclude that *Proposition A is false*. It is *not* true that the electrons go *either* through hole 1 or hole 2. But that conclusion can be tested by another experiment.

37-6 Watching the electrons

We shall now try the following experiment. To our electron apparatus we add a very strong light source, placed behind the wall and between the two holes, as shown in Fig. 37-4. We know that electric charges scatter light. So when an electron passes, however it does pass, on its way to the detector, it will scatter some light to our eye, and we can *see* where the electron goes. If, for instance, an electron were to take the path via hole 2 that is sketched in Fig. 37-4, we should see a flash of light coming from the vicinity of the place marked *A* in the figure. If an electron passes through hole 1 we would expect to see a flash from the vicinity of the upper hole. If it should happen that we get light from both places at the same time, because the electron divides in half . . . Let us just do the experiment!

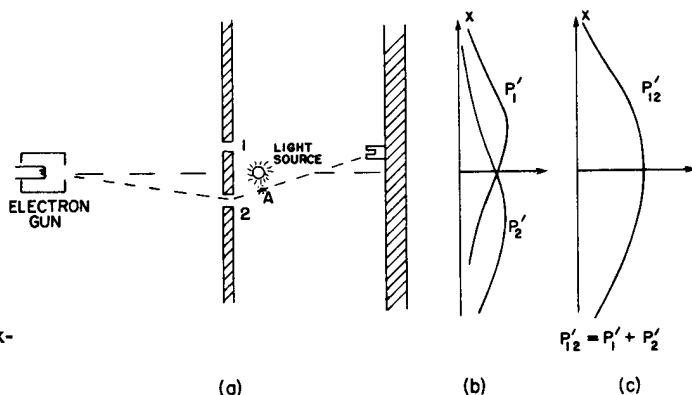


Fig. 37-4. A different electron experiment.

Here is what we see: *every* time that we hear a “click” from our electron detector (at the backstop), we *also* see a flash of light *either* near hole 1 *or* near hole 2, but *never* both at once! And we observe the same result no matter where we put the detector. From this observation we conclude that when we look at the electrons we find that the electrons go either through one hole or the other. Experimentally, Proposition A is necessarily true.

What, then, is wrong with our argument *against* Proposition A? Why *isn't* P_{12} just equal to $P_1 + P_2$? Back to experiment! Let us keep track of the electrons and find out what they are doing. For each position (x -location) of the detector we will count the electrons that arrive and *also* keep track of which hole they went through, by watching for the flashes. We can keep track of things this way: whenever we hear a “click” we will put a count in Column 1 if we see the flash near hole 1, and if we see the flash near hole 2, we will record a count in Column 2. Every electron which arrives is recorded in one of two classes: those which come through 1 and those which come through 2. From the number recorded in Column 1 we get the probability P'_1 that an electron will arrive at the detector via hole 1; and from the number recorded in Column 2 we get P'_2 , the probability that an electron will arrive at the detector via hole 2. If we now repeat such a measurement for many values of x , we get the curves for P'_1 and P'_2 shown in part (b) of Fig. 37-4.

Well, that is not too surprising! We get for P'_1 something quite similar to what we got before for P_1 by blocking off hole 2; and P'_2 is similar to what we got by blocking hole 1. So there is *not* any complicated business like going through both holes. When we watch them, the electrons come through just as we would

expect them to come through. Whether the holes are closed or open, those which we see come through hole 1 are distributed in the same way whether hole 2 is open or closed.

But wait! What do we have *now* for the *total* probability, the probability that an electron will arrive at the detector by any route? We already have that information. We just pretend that we never looked at the light flashes, and we lump together the detector clicks which we have separated into the two columns. We *must* just *add* the numbers. For the probability that an electron will arrive at the backstop by passing through *either* hole, we do find $P'_{12} = P_1 + P_2$. That is, although we succeeded in watching which hole our electrons come through, we no longer get the old interference curve P_{12} , but a new one, P'_{12} , showing no interference! If we turn out the light P_{12} is restored.

We must conclude that *when we look at the electrons* the distribution of them on the screen is different than when we do not look. Perhaps it is turning on our light source that disturbs things? It must be that the electrons are very delicate, and the light, when it scatters off the electrons, gives them a jolt that changes their motion. We know that the electric field of the light acting on a charge will exert a force on it. So perhaps we *should* expect the motion to be changed. Anyway, the light exerts a big influence on the electrons. By trying to “watch” the electrons we have changed their motions. That is, the jolt given to the electron when the photon is scattered by it is such as to change the electron’s motion enough so that if it *might* have gone to where P_{12} was at a maximum it will instead land where P_{12} was a minimum; that is why we no longer see the wavy interference effects.

You may be thinking: “Don’t use such a bright source! Turn the brightness down! The light waves will then be weaker and will not disturb the electrons so much. Surely, by making the light dimmer and dimmer, eventually the wave will be weak enough that it will have a negligible effect.” O.K. Let’s try it. The first thing we observe is that the flashes of light scattered from the electrons as they pass by does *not* get weaker. *It is always the same-sized flash.* The only thing that happens as the light is made dimmer is that sometimes we hear a “click” from the detector but see *no flash at all*. The electron has gone by without being “seen.” What we are observing is that light *also* acts like electrons, we *knew* that it was “wavy,” but now we find that it is also “lumpy.” It always arrives—or is scattered—in lumps that we call “photons.” As we turn down the *intensity* of the light source we do not change the *size* of the photons, only the *rate* at which they are emitted. *That* explains why, when our source is dim, some electrons get by without being seen. There did not happen to be a photon around at the time the electron went through.

This is all a little discouraging. If it is true that whenever we “see” the electron we see the same-sized flash, then those electrons we see are *always* the disturbed ones. Let us try the experiment with a dim light anyway. Now whenever we hear a click in the detector we will keep a count in three columns: in Column (1) those electrons seen by hole 1, in Column (2) those electrons seen by hole 2, and in Column (3) those electrons not seen at all. When we work up our data (computing the probabilities) we find these results: Those “seen by hole 1” have a distribution like P'_1 ; those “seen by hole 2” have a distribution like P'_2 (so that those “seen by either hole 1 or 2” have a distribution like P'_{12}); and those “not seen at all” have a “wavy” distribution just like P_{12} of Fig. 37-3! *If the electrons are not seen, we have interference!*

That is understandable. When we do not see the electron, no photon disturbs it, and when we do see it, a photon has disturbed it. There is always the same amount of disturbance because the light photons all produce the same-sized effects and the effect of the photons being scattered is enough to smear out any interference effect.

Is there not *some* way we can see the electrons without disturbing them? We learned in an earlier chapter that the momentum carried by a “photon” is inversely proportional to its wavelength ($p = h/\lambda$). Certainly the jolt given to the electron when the photon is scattered toward our eye depends on the momentum that photon carries. Aha! If we want to disturb the electrons only

slightly we should not have lowered the *intensity* of the light, we should have lowered its *frequency* (the same as increasing its wavelength). Let us use light of a redder color. We could even use infrared light, or radiowaves (like radar), and “see” where the electron went with the help of some equipment that can “see” light of these longer wavelengths. If we use “gentler” light perhaps we can avoid disturbing the electrons so much.

Let us try the experiment with longer waves. We shall keep repeating our experiment, each time with light of a longer wavelength. At first, nothing seems to change. The results are the same. Then a terrible thing happens. You remember that when we discussed the microscope we pointed out that, due to the *wave nature* of the light, there is a limitation on how close two spots can be and still be seen as two separate spots. This distance is of the order of the wavelength of light. So now, when we make the wavelength longer than the distance between our holes, we see a *big* fuzzy flash when the light is scattered by the electrons. We can no longer tell which hole the electron went through! We just know it went somewhere! And it is just with light of this color that we find that the jolts given to the electron are small enough so that P'_{12} begins to look like P_{12} —that we begin to get some interference effect. And it is only for wavelengths much longer than the separation of the two holes (when we have no chance at all of telling where the electron went) that the disturbance due to the light gets sufficiently small that we again get the curve P_{12} shown in Fig. 37-3.

In our experiment we find that it is impossible to arrange the light in such a way that one can tell which hole the electron went through, and at the same time not disturb the pattern. It was suggested by Heisenberg that the then new laws of nature could only be consistent if there were some basic limitation on our experimental capabilities not previously recognized. He proposed, as a general principle, his *uncertainty principle*, which we can state in terms of our experiment as follows: “It is impossible to design an apparatus to determine which hole the electron passes through, that will not at the same time disturb the electrons enough to destroy the interference pattern.” If an apparatus is capable of determining which hole the electron goes through, it *cannot* be so delicate that it does not disturb the pattern in an essential way. No one has ever found (or even thought of) a way around the uncertainty principle. So we must assume that it describes a basic characteristic of nature.

The complete theory of quantum mechanics which we now use to describe atoms and, in fact, all matter depends on the correctness of the uncertainty principle. Since quantum mechanics is such a successful theory, our belief in the uncertainty principle is reinforced. But if a way to “beat” the uncertainty principle were ever discovered, quantum mechanics would give inconsistent results and would have to be discarded as a valid theory of nature.

“Well,” you say, “what about Proposition A? It is true, or is it *not* true, that the electron either goes through hole 1 or it goes through hole 2?” The only answer that can be given is that we have found from experiment that there is a certain special way that we have to think in order that we do not get into inconsistencies. What we must say (to avoid making wrong predictions) is the following. If one looks at the holes or, more accurately, if one has a piece of apparatus which is capable of determining whether the electrons go through hole 1 or hole 2, then one *can* say that it goes either through hole 1 or hole 2. *But*, when one does *not* try to tell which way the electron goes, when there is nothing in the experiment to disturb the electrons, then one may *not* say that an electron goes either through hole 1 or hole 2. If one does say that, and starts to make any deductions from the statement, he will make errors in the analysis. This is the logical tightrope on which we must walk if we wish to describe nature successfully.

If the motion of all matter—as well as electrons—must be described in terms of waves, what about the bullets in our first experiment? Why didn’t we see an interference pattern there? It turns out that for the bullets the wavelengths were so tiny that the interference patterns became very fine. So fine, in fact, that with any

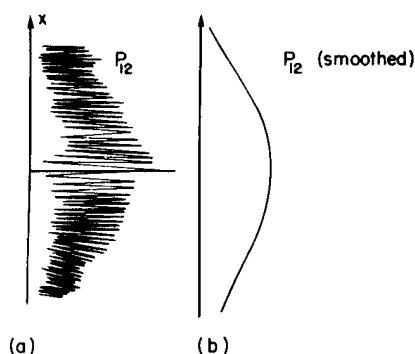


Fig. 37-5. Interference pattern with bullets: (a) actual (schematic), (b) observed.

detector of finite size one could not distinguish the separate maxima and minima. What we saw was only a kind of average, which is the classical curve. In Fig. 37-5 we have tried to indicate schematically what happens with large-scale objects. Part (a) of the figure shows the probability distribution one might predict for bullets, using quantum mechanics. The rapid wiggles are supposed to represent the interference pattern one gets for waves of very short wavelength. Any physical detector, however, straddles several wiggles of the probability curve, so that the measurements show the smooth curve drawn in part (b) of the figure.

37-7 First principles of quantum mechanics

We will now write a summary of the main conclusions of our experiments. We will, however, put the results in a form which makes them true for a general class of such experiments. We can write our summary more simply if we first define an “ideal experiment” as one in which there are no uncertain external influences, i.e., no jiggling or other things going on that we cannot take into account. We would be quite precise if we said: “An ideal experiment is one in which all of the initial and final conditions of the experiment are completely specified.” What we will call “an event” is, in general, just a specific set of initial and final conditions. (For example: “an electron leaves the gun, arrives at the detector, and nothing else happens.”) Now for our summary.

SUMMARY

- (1) The probability of an event in an ideal experiment is given by the square of the absolute value of a complex number ϕ which is called the probability amplitude.

$$\begin{aligned} P &= \text{probability,} \\ \phi &= \text{probability amplitude,} \\ P &= |\phi|^2. \end{aligned} \tag{37.6}$$

- (2) When an event can occur in several alternative ways, the probability amplitude for the event is the sum of the probability amplitudes for each way considered separately. There is interference.

$$\begin{aligned} \phi &= \phi_1 + \phi_2, \\ P &= |\phi_1 + \phi_2|^2. \end{aligned} \tag{37.7}$$

- (3) If an experiment is performed which is capable of determining whether one or another alternative is actually taken, the probability of the event is the sum of the probabilities for each alternative. The interference is lost.

$$P = P_1 + P_2. \tag{37.8}$$

One might still like to ask: “How does it work? What is the machinery behind the law?” No one has found any machinery behind the law. No one can “explain” any more than we have just “explained.” No one will give you any deeper representation of the situation. We have no ideas about a more basic mechanism from which these results can be deduced.

We would like to emphasize a very important difference between classical and quantum mechanics. We have been talking about the probability that an electron will arrive in a given circumstance. We have implied that in our experimental arrangement (or even in the best possible one) it would be impossible to predict exactly what would happen. We can only predict the odds! This would mean, if it were true, that physics has given up on the problem of trying to predict exactly what will happen in a definite circumstance. Yes! physics *has* given up. *We do not know how to predict what would happen in a given circumstance*, and we believe now that it is impossible, that the only thing that can be predicted is the probability of different events. It must be recognized that this is a retrenchment in our earlier ideal of understanding nature. It may be a backward step, but no one has seen a way to avoid it.

We make now a few remarks on a suggestion that has sometimes been made to try to avoid the description we have given: "Perhaps the electron has some kind of internal works—some inner variables—that we do not yet know about. Perhaps that is why we cannot predict what will happen. If we could look more closely at the electron we could be able to tell where it would end up." So far as we know, that is impossible. We would still be in difficulty. Suppose we were to assume that inside the electron there is some kind of machinery that determines where it is going to end up. That machine must *also* determine which hole it is going to go through on its way. But we must not forget that what is inside the electron should not be dependent on what *we* do, and in particular upon whether we open or close one of the holes. So if an electron, before it starts, has already made up its mind (a) which hole it is going to use, and (b) where it is going to land, we should find P_1 for those electrons that have chosen hole 1, P_2 for those that have chosen hole 2, and *necessarily* the sum $P_1 + P_2$ for those that arrive through the two holes. There seems to be no way around this. But we have verified experimentally that that is not the case. And no one has figured a way out of this puzzle. So at the present time we must limit ourselves to computing probabilities. We say "at the present time," but we suspect very strongly that it is something that will be with us forever—that it is impossible to beat that puzzle—that this is the way nature really *is*.

37-8 The uncertainty principle

This is the way Heisenberg stated the uncertainty principle originally: If you make the measurement on any object, and you can determine the x -component of its momentum with an uncertainty Δp , you cannot, at the same time, know its x -position more accurately than $\Delta x = h/\Delta p$. The uncertainties in the position and momentum at any instant must have their product greater than Planck's constant. This is a special case of the uncertainty principle that was stated above more generally. The more general statement was that one cannot design equipment in any way to determine which of two alternatives is taken, without, at the same time, destroying the pattern of interference.

Let us show for one particular case that the kind of relation given by Heisenberg must be true in order to keep from getting into trouble. We imagine a modification of the experiment of Fig. 37-3, in which the wall with the holes consists of a plate mounted on rollers so that it can move freely up and down (in the x -direction), as shown in Fig. 37-6. By watching the motion of the plate carefully we can try to tell which hole an electron goes through. Imagine what happens when the detector is placed at $x = 0$. We would expect that an electron which passes through hole 1 must be deflected downward by the plate to reach the detector. Since the vertical component of the electron momentum is changed, the plate must recoil with an equal momentum in the opposite direction. The plate will get an upward kick. If the electron goes through the lower hole, the plate should feel a downward kick. It is clear that for every position of the detector, the momentum received by the plate will have a different value for a traversal via hole 1 than for a traversal via hole 2. So! Without disturbing the electrons *at all*, but just by watching the *plate*, we can tell which path the electron used.

Now in order to do this it is necessary to know what the momentum of the screen is, before the electron goes through. So when we measure the momentum after the electron goes by, we can figure out how much the plate's momentum has changed. But remember, according to the uncertainty principle we cannot at the same time know the position of the plate with an arbitrary accuracy. But if we do not know exactly *where* the plate is we cannot say precisely where the two holes are. They will be in a different place for every electron that goes through. This means that the center of our interference pattern will have a different location for each electron. The wiggles of the interference pattern will be smeared out. We shall show quantitatively in the next chapter that if we determine the momentum of the plate sufficiently accurately to determine from the recoil measurement which hole was used, then the uncertainty in the x -position of the plate will, according to the un-

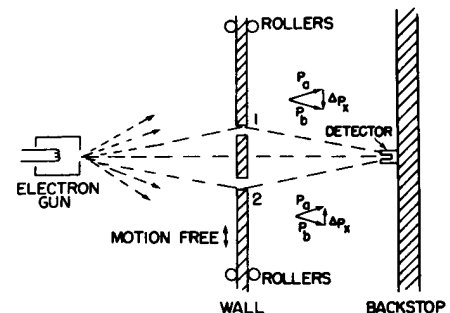


Fig. 37-6. An experiment in which the recoil of the wall is measured.

certainty principle, be enough to shift the pattern observed at the detector up and down in the x -direction about the distance from a maximum to its nearest minimum. Such a random shift is just enough to smear out the pattern so that no interference is observed.

The uncertainty principle “protects” quantum mechanics. Heisenberg recognized that if it were possible to measure the momentum and the position simultaneously with a greater accuracy, the quantum mechanics would collapse. So he proposed that it must be impossible. Then people sat down and tried to figure out ways of doing it, and nobody could figure out a way to measure the position and the momentum of anything—a screen, an electron, a billiard ball, anything—with any greater accuracy. Quantum mechanics maintains its perilous but accurate existence.

The Relation of Wave and Particle Viewpoints

38-1 Probability wave amplitudes

In this chapter we shall discuss the relationship of the wave and particle viewpoints. We already know, from the last chapter, that neither the wave viewpoint nor the particle viewpoint is correct. Usually we have tried to present things accurately, or at least precisely enough that they will not have to be changed when we learn more—it may be extended, but it will not be changed! But when we try to talk about the wave picture or the particle picture, both are approximate, and both will change. Therefore what we learn in this chapter will not be accurate in a certain sense; it is a kind of half-intuitive argument that will be made more precise later, but certain things will be changed a little bit when we interpret them correctly in quantum mechanics. The reason for doing such a thing, of course, is that we are not going to go directly into quantum mechanics, but we want to have at least some idea of the kinds of effects that we will find. Furthermore, all our experiences are with waves and with particles, and so it is rather handy to use the wave and particle ideas to get some understanding of what happens in given circumstances before we know the complete mathematics of the quantum-mechanical amplitudes. We shall try to illustrate the weakest places as we go along, but most of it is very nearly correct—it is just a matter of interpretation.

First of all, we know that the new way of representing the world in quantum mechanics—the new framework—is to give an amplitude for every event that can occur, and if the event involves the reception of one particle then we can give the amplitude to find that one particle at different places and at different times. The probability of finding the particle is then proportional to the absolute square of the amplitude. In general, the amplitude to find a particle in different places at different times varies with position and time.

In a special case the amplitude varies sinusoidally in space and time like $e^{i(\omega t - \mathbf{k} \cdot \mathbf{r})}$ (do not forget that these amplitudes are complex numbers, not real numbers) and involves a definite frequency ω and wave number \mathbf{k} . Then it turns out that this corresponds to a classical limiting situation where we would have believed that we have a particle whose energy E was known and is related to the frequency by

$$E = \hbar\omega, \quad (38.1)$$

and whose momentum \mathbf{p} is also known and is related to the wave number by

$$\mathbf{p} = \hbar\mathbf{k}. \quad (38.2)$$

This means that the idea of a particle is limited. The idea of a particle—its location, its momentum, etc.—which we use so much, is in certain ways unsatisfactory. For instance, if an amplitude to find a particle at different places is given by $e^{i(\omega t - \mathbf{k} \cdot \mathbf{r})}$, whose absolute square is a constant, that would mean that the probability of finding a particle is the same at all points. That means we do not know *where* it is—it can be anywhere—there is a great uncertainty in its location.

On the other hand, if the position of a particle is more or less well known and we can predict it fairly accurately, then the probability of finding it in different places must be confined to a certain region, whose length we call Δx . Outside this region, the probability is zero. Now this probability is the absolute square of an amplitude, and if the absolute square is zero, the amplitude is also zero, so that

38-1 Probability wave amplitudes

38-2 Measurement of position and momentum

38-3 Crystal diffraction

38-4 The size of an atom

38-5 Energy levels

38-6 Philosophical implications

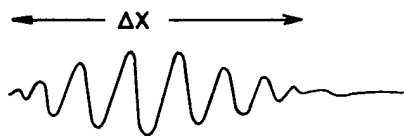


Fig. 38-1. A wave packet of length Δx .

we have a wave train whose length is Δx (Fig. 38-1), and the wavelength (the distance between nodes of the waves in the train) of that wave train is what corresponds to the particle momentum.

Here we encounter a strange thing about waves; a very simple thing which has nothing to do with quantum mechanics strictly. It is something that anybody who works with waves, even if he knows no quantum mechanics, knows: namely, *we cannot define a unique wavelength for a short wave train*. Such a wave train does not have a definite wavelength; there is an indefiniteness in the wave number that is related to the finite length of the train, and thus there is an indefiniteness in the momentum.

38-2 Measurement of position and momentum

Let us consider two examples of this idea—to see the reason why there is an uncertainty in the position and/or the momentum, if quantum mechanics is right. We have also seen before that if there were not such a thing—if it were possible to measure the position and the momentum of anything simultaneously—we would have a paradox; it is fortunate that we do not have such a paradox, and the fact that such an uncertainty comes naturally from the wave picture shows that everything is mutually consistent.

Here is one example which shows the relationship between the position and the momentum in a circumstance that is easy to understand. Suppose we have a single slit, and particles are coming from very far away with a certain energy—so that they are all coming essentially horizontally (Fig. 38-2). We are going to concentrate on the vertical components of momentum. All of these particles have a certain horizontal momentum p_0 , say, in a classical sense. So, in the classical sense, the vertical momentum p_y , before the particle goes through the hole, is definitely known. The particle is moving neither up nor down, because it came from a source that is far away—and so the vertical momentum is of course zero. But now let us suppose that it goes through a hole whose width is B . Then after it has come out through the hole, we know the position vertically—the y position—with considerable accuracy—namely $\pm B$. That is, the uncertainty in position, Δy , is of order B . Now we might also want to say, since we know the momentum is absolutely horizontal, that Δp_y is zero; but that is wrong. We *once* knew the momentum was horizontal, but we do not know it any more. Before the particles passed through the hole, we did not know their vertical positions. Now that we have found the vertical position by having the particle come through the hole, we have lost our information on the vertical momentum! Why? According to the wave theory, there is a spreading out, or diffraction, of the waves after they go through the slit, just as for light. Therefore there is a certain probability that particles coming out of the slit are not coming exactly straight. The pattern is spread out by the diffraction effect, and the angle of spread, which we can define as the angle of the first minimum, is a measure of the uncertainty in the final angle.

How does the pattern become spread? To say it is spread means that there is some chance for the particle to be moving up or down, that is, to have a component of momentum up or down. We say *chance* and *particle* because we can detect this diffraction pattern with a particle counter, and when the counter receives the particle, say at C in Fig. 38-2, it receives the *entire* particle, so that, in a classical sense, the particle has a vertical momentum, in order to get from the slit up to C .

To get a rough idea of the spread of the momentum, the vertical momentum p_y has a spread which is equal to $p_0 \Delta\theta$, where p_0 is the horizontal momentum. And how big is $\Delta\theta$ in the spread-out pattern? We know that the first minimum occurs at an angle $\Delta\theta$ such that the waves from one edge of the slit have to travel one wavelength farther than the waves from the other side—we worked that out before (Chapter 30). Therefore $\Delta\theta$ is λ/B , and so Δp_y in this experiment is $p_0 \lambda/B$. Note that if we make B smaller and make a more accurate measurement of the position of the particle, the diffraction pattern gets wider. Remember, when we closed the slits on the experiment with the microwaves, we had more intensity farther out. So the narrower we make the slit, the wider the pattern gets, and the

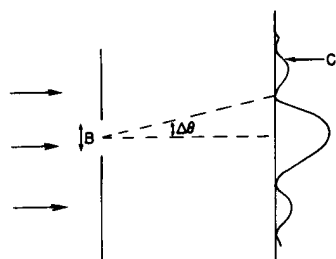


Fig. 38-2. Diffraction of particles passing through a slit.

more is the likelihood that we would find that the particle has sidewise momentum. Thus the uncertainty in the vertical momentum is inversely proportional to the uncertainty of y . In fact, we see that the product of the two is equal to $p_0\lambda$. But λ is the wavelength and p_0 is the momentum, and in accordance with quantum mechanics, the wavelength times the momentum is Planck's constant h . So we obtain the rule that the uncertainties in the vertical momentum and in the vertical position have a product of the order h :

$$\Delta y \Delta p_y \approx h. \quad (38.3)$$

We cannot prepare a system in which we know the vertical position of a particle and can predict how it will move vertically with greater certainty than given by (38.3). That is, the uncertainty in the vertical momentum must exceed $h/\Delta y$, where Δy is the uncertainty in our knowledge of the position.

Sometimes people say quantum mechanics is all wrong. When the particle arrived from the left, its vertical momentum was zero. And now that it has gone through the slit, its position is known. Both position and momentum seem to be known with arbitrary accuracy. It is quite true that we can receive a particle, and on reception determine what its position is and what its momentum would have had to have been to have gotten there. That is true, but that is not what the uncertainty relation (38.3) refers to. Equation (38.3) refers to the *predictability* of a situation, not remarks about the *past*. It does no good to say "I knew what the momentum was before it went through the slit, and now I know the position," because now the momentum knowledge is lost. The fact that it went through the slit no longer permits us to predict the vertical momentum. We are talking about a predictive theory, not just measurements after the fact. So we must talk about what we can predict.

Now let us take the thing the other way around. Let us take another example of the same phenomenon, a little more quantitatively. In the previous example we measured the momentum by a classical method. Namely, we considered the direction and the velocity and the angles, etc., so we got the momentum by classical analysis. But since momentum is related to wave number, there exists in nature still another way to measure the momentum of a particle—photon or otherwise—which has no classical analog, because it uses Eq. (38.2). We measure the *wavelengths of the waves*. Let us try to measure momentum in this way.

Suppose we have a grating with a large number of lines (Fig. 38-3), and send a beam of particles at the grating. We have often discussed this problem: if the particles have a definite momentum, then we get a very sharp pattern in a certain direction, because of the interference. And we have also talked about how accurately we can determine that momentum, that is to say, what the resolving power of such a grating is. Rather than derive it again, we refer to Chapter 30, where we found that the relative uncertainty in the wavelength that can be measured with a given grating is $1/Nm$, where N is the number of lines on the grating and m is the order of the diffraction pattern. That is,

$$\Delta\lambda/\lambda = 1/Nm. \quad (38.4)$$

Now formula (38.4) can be rewritten as

$$\Delta\lambda/\lambda^2 = 1/Nm\lambda = 1/L, \quad (38.5)$$

where L is the distance shown in Fig. 38-3. This distance is the difference between the total distance that the particle or wave or whatever it is has to travel if it is reflected from the bottom of the grating, and the distance that it has to travel if it is reflected from the top of the grating. That is, the waves which form the diffraction pattern are waves which come from different parts of the grating. The first ones that arrive come from the bottom end of the grating, from the beginning of the wave train, and the rest of them come from later parts of the wave train, coming from different parts of the grating, until the last one finally arrives, and that involves a point in the wave train a distance L behind the first point. So in order that we shall have a sharp line in our spectrum corresponding to a definite momentum,

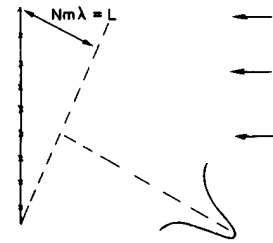


Fig. 38-3. Determination of momentum by using a diffraction grating.

with an uncertainty given by (38.4), we have to have a wave train of at least length L . If the wave train is too short we are not using the entire grating. The waves which form the spectrum are being reflected from only a very short sector of the grating if the wave train is too short, and the grating will not work right—we will find a big angular spread. In order to get a narrower one, we need to use the whole grating, so that at least at some moment the whole wave train is scattering simultaneously from all parts of the grating. Thus the wave train must be of length L in order to have an uncertainty in the wavelength less than that given by (38.5). Incidentally,

$$\Delta\lambda/\lambda^2 = \Delta(1/\lambda) = \Delta k/2\pi. \quad (38.6)$$

Therefore

$$\Delta k = 2\pi/L, \quad (38.7)$$

where L is the length of the wave train.

This means that if we have a wave train whose length is less than L , the uncertainty in the wave number must exceed $2\pi/L$. Or the uncertainty in a wave number times the length of the wave train—we will call that for a moment Δx —exceeds 2π . We call it Δx because that is the uncertainty in the location of the particle. If the wave train exists only in a finite length, then that is where we could find the particle, within an uncertainty Δx . Now this property of waves, that the length of the wave train times the uncertainty of the wave number associated with it is at least 2π , is a property that is known to everyone who studies them. It has nothing to do with quantum mechanics. It is simply that if we have a finite train, we cannot count the waves in it very precisely. Let us try another way to see the reason for that.

Suppose that we have a finite train of length L ; then because of the way it has to decrease at the ends, as in Fig. 38-1, the number of waves in the length L is uncertain by something like ± 1 . But the number of waves in L is $kL/2\pi$. Thus k is uncertain, and we again get the result (38.7), a property merely of waves. The same thing works whether the waves are in space and k is the number of radians per centimeter and L is the length of the train, or the waves are in time and ω is the number of oscillations per second and T is the “length” in time that the wave train comes in. That is, if we have a wave train lasting only for a certain finite time T , then the uncertainty in the frequency is given by

$$\Delta\omega = 2\pi/T. \quad (38.8)$$

We have tried to emphasize that these are properties of waves alone, and they are well known, for example, in the theory of sound.

The point is that in quantum mechanics we interpret the wave number as being a measure of the momentum of a particle, with the rule that $p = \hbar k$, so that relation (38.7) tells us that $\Delta p \approx \hbar/\Delta x$. This, then, is a limitation of the classical idea of momentum. (Naturally, it has to be limited in some ways if we are going to represent particles by waves!) It is nice that we have found a rule that gives us some idea of when there is a failure of classical ideas.

38-3 Crystal diffraction

Next let us consider the reflection of particle waves from a crystal. A crystal is a thick thing which has a whole lot of similar atoms—we will include some complications later—in a nice array. The question is how to set the array so that we get a strong reflected maximum in a given direction for a given beam of, say, light (x-rays), electrons, neutrons, or anything else. In order to obtain a strong reflection, the scattering from all of the atoms must be in phase. There cannot be equal numbers in phase and out of phase, or the waves will cancel out. The way to arrange things is to find the regions of constant phase, as we have already explained; they are planes which make equal angles with the initial and final directions (Fig. 38-4).

If we consider two parallel planes, as in Fig. 38-4, the waves scattered from the two planes will be in phase provided the difference in distance travelled by a wave-

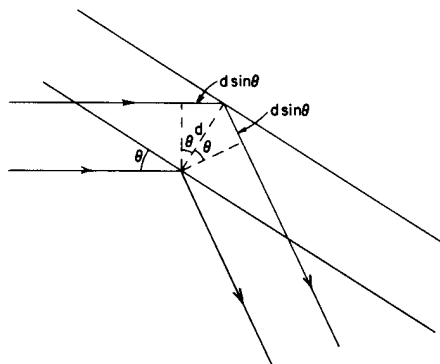


Fig. 38-4. Scattering of waves by crystal planes.

front is an integral number of wavelengths. This difference can be seen to be $2d \sin \theta$, where d is the perpendicular distance between the planes. Thus the condition for coherent reflection is

$$2d \sin \theta = n\lambda \quad (n = 1, 2, \dots). \quad (38.9)$$

If, for example, the crystal is such that the atoms happen to lie on planes obeying condition (38.9) with $n = 1$, then there will be a strong reflection. If, on the other hand, there are other atoms of the same nature (equal in density) halfway between, then the intermediate planes will also scatter equally strongly and will interfere with the others and produce no effect. So d in (38.9) must refer to *adjacent* planes; we cannot take a plane five layers farther back and use this formula!

As a matter of interest, actual crystals are not usually as simple as a single kind of atom repeated in a certain way. Instead, if we make a two-dimensional analog, they are much like wallpaper, in which there is some kind of figure which repeats all over the wallpaper. By "figure" we mean, in the case of atoms, some arrangement—calcium and a carbon and three oxygens, etc., for calcium carbonate, and so on—which may involve a relatively large number of atoms. But whatever it is, the figure is repeated in a pattern. This basic figure is called a *unit cell*.

The basic pattern of repetition defines what we call the *lattice type*; the lattice type can be immediately determined by looking at the reflections and seeing what their symmetry is. In other words, where we find any reflections *at all* determines the lattice type, but in order to determine what is in each of the elements of the lattice one must take into account the *intensity* of the scattering at the various directions. *Which* directions scatter depends on the type of lattice, but *how strongly* each scatters is determined by what is inside each unit cell, and in that way the structure of crystals is worked out.

Two photographs of x-ray diffraction patterns are shown in Figs. 38-5 and 38-6; they illustrate scattering from rock salt and myoglobin, respectively.

Incidentally, an interesting thing happens if the spacings of the nearest planes are less than $\lambda/2$. In this case (38.9) has no solution for n . Thus if λ is bigger than twice the distance between adjacent planes then there is no side diffraction pattern, and the light—or whatever it is—will go right through the material without bouncing off or getting lost. So in the case of light, where λ is much bigger than the spacing, of course it does go through and there is no pattern of reflection from the planes of the crystal.

This fact also has an interesting consequence in the case of piles which make neutrons (these are obviously particles, for anybody's money!). If we take these neutrons and let them into a long block of graphite, the neutrons diffuse and work their way along (Fig. 38-7). They diffuse because they are bounced by the atoms, but strictly, in the wave theory, they are bounced by the atoms because of diffraction from the crystal planes. It turns out that if we take a very long piece of graphite, the neutrons that come out the far end are all of long wavelength! In fact, if one plots the intensity as a function of wavelength, we get nothing except for wavelengths longer than a certain minimum (Fig. 38-8). In other words, we can get very slow neutrons that way. Only the slowest neutrons come through; they are not diffracted or scattered by the crystal planes of the graphite, but keep going right through like light through glass, and are not scattered out the sides. There are many other demonstrations of the reality of neutron waves and waves of other particles.

38-4 The size of an atom

We now consider another application of the uncertainty relation, Eq. (38.3). It must not be taken too seriously; the idea is right but the analysis is not very accurate. The idea has to do with the determination of the size of atoms, and the fact that, classically, the electrons would radiate light and spiral in until they settle down right on top of the nucleus. But that cannot be right quantum-mechanically because then we would know where each electron was and how fast it was moving.

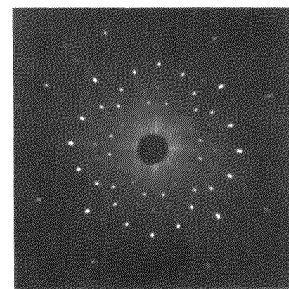


Figure 38-5

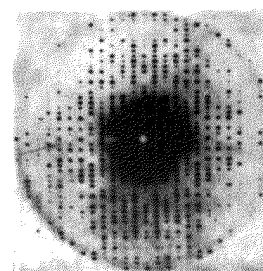


Figure 38-6

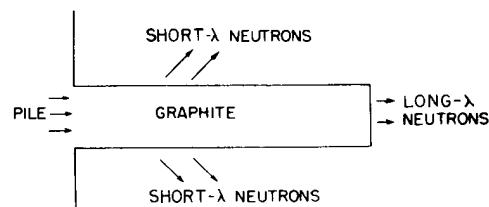


Fig. 38-7. Diffusion of pile neutrons through graphite block.

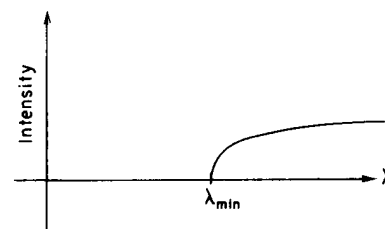


Fig. 38-8. Intensity of neutrons out of graphite rod as function of wavelength.

Suppose we have a hydrogen atom, and measure the position of the electron; we must not be able to predict exactly where the electron will be, or the momentum spread will then turn out to be infinite. Every time we look at the electron, it is somewhere, but it has an amplitude to be in different places so there is a probability of it being found in different places. These places cannot all be at the nucleus; we shall suppose there is a spread in position of order a . That is, the distance of the electron from the nucleus is usually about a . We shall determine a by minimizing the total energy of the atom.

The spread in momentum is roughly h/a because of the uncertainty relation, so that if we try to measure the momentum of the electron in some manner, such as by scattering x-rays off it and looking for the Doppler effect from a moving scatterer, we would expect not to get zero every time—the electron is not standing still—but the momenta must be of the order $p \approx h/a$. Then the kinetic energy is roughly $\frac{1}{2}mv^2 = p^2/2m = h^2/2ma^2$. (In a sense, this is a kind of dimensional analysis to find out in what way the kinetic energy depends upon Planck's constant, upon m , and upon the size of the atom. We need not trust our answer to within factors like 2, π , etc. We have not even defined a very precisely.) Now the potential energy is minus e^2 over the distance from the center, say $-e^2/a$, where, we remember, e^2 is the charge of an electron squared, divided by $4\pi\epsilon_0$. Now the point is that the potential energy is reduced if a gets smaller, but the smaller a is, the higher the momentum required, because of the uncertainty principle, and therefore the higher the kinetic energy. The total energy is

$$E = h^2/2ma^2 - e^2/a. \quad (38.10)$$

We do not know what a is, but we know that the atom is going to arrange itself to make some kind of compromise so that the energy is as little as possible. In order to minimize E , we differentiate with respect to a , set the derivative equal to zero, and solve for a . The derivative of E is

$$dE/da = -h^2/ma^3 + e^2/a^2, \quad (38.11)$$

and setting $dE/da = 0$ gives for a the value

$$\begin{aligned} a_0 &= h^2/me^2 = 0.528 \text{ angstrom} \\ &= 0.528 \times 10^{-10} \text{ meter.} \end{aligned} \quad (38.12)$$

This particular distance is called the *Bohr radius*, and we have thus learned that atomic dimensions are of the order of angstroms, which is right: This is pretty good—in fact, it is amazing, since until now we have had no basis for understanding the size of atoms! Atoms are completely impossible from the classical point of view, since the electrons would spiral into the nucleus.

Now if we put the value (38.12) for a_0 into (38.10) to find the energy, it comes out

$$E_0 = -e^2/2a_0 = -me^4/2h^2 = -13.6 \text{ ev.} \quad (38.13)$$

What does a negative energy mean? It means that the electron has less energy when it is in the atom than when it is free. It means it is bound. It means it takes energy to kick the electron out; it takes energy of the order of 13.6 ev to ionize a hydrogen atom. We have no reason to think that it is not two or three times this—or half of this—or $(1/\pi)$ times this, because we have used such a sloppy argument. However, we have cheated, we have used all the constants in such a way that it happens to come out the right number! This number, 13.6 electron volts, is called a Rydberg of energy; it is the ionization energy of hydrogen.

So we now understand why we do not fall through the floor. As we walk, our shoes with their masses of atoms push against the floor with *its* mass of atoms. In order to squash the atoms closer together, the electrons would be confined to a smaller space and, by the uncertainty principle, their momenta would have to be higher on the average, and that means high energy; the resistance to atomic compression is a quantum-mechanical effect and not a classical effect. Classically,

we would expect that if we were to draw all the electrons and protons closer together, the energy would be reduced still further, and the best arrangement of positive and negative charges in classical physics is all on top of each other. This was well known in classical physics and was a puzzle because of the existence of the atom. Of course, the early scientists invented some ways out of the trouble—but never mind, we have the *right* way out, now! (Maybe.)

Incidentally, although we have no reason to understand it at the moment, in a situation where there are many electrons it turns out that they try to keep away from each other. If one electron is occupying a certain space, then another does not occupy the same space. More precisely, there are two spin cases, so that two can sit on top of each other, one spinning one way and one the other way. But after that we cannot put any more there. We have to put others in another place, and that is the real reason that matter has strength. If we could put all the electrons in the same place it would condense even more than it does. It is the fact that the electrons cannot all get on top of each other that makes tables and everything else solid.

Obviously, in order to understand the properties of matter, we will have to use quantum mechanics and not be satisfied with classical mechanics.

38-5 Energy levels

We have talked about the atom in its lowest possible energy condition, but it turns out that the electron can do other things. It can jiggle and wiggle in a more energetic manner, and so there are many different possible motions for the atom. According to quantum mechanics, in a stationary condition there can only be definite energies for an atom. We make a diagram (Fig. 38-9) in which we plot the energy vertically, and we make a horizontal line for each allowed value of the energy. When the electron is free, i.e., when its energy is positive, it can have any energy; it can be moving at any speed. But bound energies are not arbitrary. The atom must have one or another out of a set of allowed values, such as those in Fig. 38-9.

Now let us call the allowed values of the energy E_0, E_1, E_2, E_3 . If an atom is initially in one of these “excited states,” E_1, E_2 , etc., it does not remain in that state forever. Sooner or later it drops to a lower state and radiates energy in the form of light. The frequency of the light that is emitted is determined by conservation of energy plus the quantum-mechanical understanding that the frequency of the light is related to the energy of the light by (38.1). Therefore the frequency of the light which is liberated in a transition from energy E_3 to energy E_1 (for example) is

$$\omega_{31} = (E_3 - E_1)/\hbar. \quad (38.14)$$

This, then, is a characteristic frequency of the atom and defines a spectral emission line. Another possible transition would be from E_3 to E_0 . That would have a different frequency

$$\omega_{30} = (E_3 - E_0)/\hbar. \quad (38.15)$$

Another possibility is that if the atom were excited to the state E_1 it could drop to the ground state E_0 , emitting a photon of frequency

$$\omega_{10} = (E_1 - E_0)/\hbar. \quad (38.16)$$

The reason we bring up three transitions is to point out an interesting relationship. It is easy to see from (38.14), (38.15), and (38.16) that

$$\omega_{30} = \omega_{31} + \omega_{10}. \quad (38.17)$$

In general, if we find two spectral lines, we shall expect to find another line at the sum of the frequencies (or the difference in the frequencies), and that all the lines can be understood by finding a series of levels such that every line corresponds to the difference in energy of some pair of levels. This remarkable coincidence in

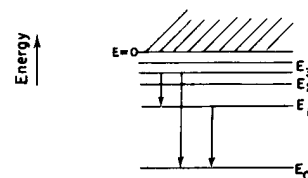


Fig. 38-9. Energy diagram for an atom, showing several possible transitions.

spectral frequencies was noted before quantum mechanics was discovered, and it is called the *Ritz combination principle*. This is again a mystery from the point of view of classical mechanics. Let us not belabor the point that classical mechanics is a failure in the atomic domain; we seem to have demonstrated that pretty well.

We have already talked about quantum mechanics as being represented by amplitudes which behave like waves, with certain frequencies and wave numbers. Let us observe how it comes about from the point of view of amplitudes that the atom has definite energy states. This is something we cannot understand from what has been said so far, but we are all familiar with the fact that confined waves have definite frequencies. For instance, if sound is confined to an organ pipe, or anything like that, then there is more than one way that the sound can vibrate, but for each such way there is a definite frequency. Thus an object in which the waves are confined has certain resonance frequencies. It is therefore a property of waves in a confined space—a subject which we will discuss in detail with formulas later on—that they exist only at definite frequencies. And since the general relation exists between frequencies of the amplitude and energy, we are not surprised to find definite energies associated with electrons bound in atoms.

38-6 Philosophical implications

Let us consider briefly some philosophical implications of quantum mechanics. As always, there are two aspects of the problem: one is the philosophical implication for physics, and the other is the extrapolation of philosophical matters to other fields. When philosophical ideas associated with science are dragged into another field, they are usually completely distorted. Therefore we shall confine our remarks as much as possible to physics itself.

First of all, the most interesting aspect is the idea of the uncertainty principle; making an observation affects the phenomenon. It has always been known that making observations affects a phenomenon, but the point is that the effect cannot be disregarded or minimized or decreased arbitrarily by rearranging the apparatus. When we look for a certain phenomenon we cannot help but disturb it in a certain minimum way, and *the disturbance is necessary for the consistency of the viewpoint*. The observer was sometimes important in prequantum physics, but only in a rather trivial sense. The problem has been raised: if a tree falls in a forest and there is nobody there to hear it, does it make a noise? A *real* tree falling in a *real* forest makes a sound, of course, even if nobody is there. Even if no one is present to hear it, there are other traces left. The sound will shake some leaves, and if we were careful enough we might find somewhere that some thorn had rubbed against a leaf and made a tiny scratch that could not be explained unless we assumed the leaf were vibrating. So in a certain sense we would have to admit that there is a sound made. We might ask: was there a *sensation* of sound? No, sensations have to do, presumably, with consciousness. And whether ants are conscious and whether there were ants in the forest, or whether the tree was conscious, we do not know. Let us leave the problem in that form.

Another thing that people have emphasized since quantum mechanics was developed is the idea that we should not speak about those things which we cannot measure. (Actually relativity theory also said this.) Unless a thing can be defined by measurement, it has no place in a theory. And since an accurate value of the momentum of a localized particle cannot be defined by measurement it therefore has no place in the theory. The idea that this is what was the matter with classical theory is a *false position*. It is a careless analysis of the situation. Just because we cannot *measure* position and momentum precisely does not *a priori* mean that we *cannot* talk about them. It only means that we *need* not talk about them. The situation in the sciences is this: A concept or an idea which cannot be measured or cannot be referred directly to experiment may or may not be useful. It need not exist in a theory. In other words, suppose we compare the classical theory of the world with the quantum theory of the world, and suppose that it is true experimentally that we can measure position and momentum only imprecisely. The question is whether the *ideas* of the exact position of a particle and the exact

momentum of a particle are valid or not. The classical theory admits the ideas; the quantum theory does not. This does not in itself mean that classical physics is wrong. When the new quantum mechanics was discovered, the classical people—which included everybody except Heisenberg, Schrödinger, and Born—said: “Look, your theory is not any good because you cannot answer certain questions like: what is the exact position of a particle?, which hole does it go through?, and some others.” Heisenberg’s answer was: “I do not need to answer such questions because you cannot ask such a question experimentally.” It is that we do not *have* to. Consider two theories (a) and (b); (a) contains an idea that cannot be checked directly but which is used in the analysis, and the other, (b), does not contain the idea. If they disagree in their predictions, one could not claim that (b) is false because it cannot explain this idea that is in (a), because that idea is one of the things that cannot be checked directly. It is always good to know which ideas cannot be checked directly, but it is not necessary to remove them all. It is not true that we can pursue science completely by using only those concepts which are directly subject to experiment.

In quantum mechanics itself there is a wave function amplitude, there is a potential, and there are many constructs that we cannot measure directly. The basis of a science is its ability to *predict*. To predict means to tell what will happen in an experiment that has never been done. How can we do that? By assuming that we know what is there, independent of the experiment. We must extrapolate the experiments to a region where they have not been done. We must take our concepts and extend them to places where they have not yet been checked. If we do not do that, we have no prediction. So it was perfectly sensible for the classical physicists to go happily along and suppose that the position—which obviously means something for a baseball—meant something also for an electron. It was not stupidity. It was a sensible procedure. Today we say that the law of relativity is supposed to be true at all energies, but someday somebody may come along and say how stupid we were. We do not know where we are “stupid” until we “stick our neck out,” and so the whole idea is to put our neck out. And the only way to find out that we are wrong is to find out *what* our predictions are. It is absolutely necessary to make constructs.

We have already made a few remarks about the indeterminacy of quantum mechanics. That is, that we are unable now to predict what will happen in physics in a given physical circumstance which is arranged as carefully as possible. If we have an atom that is in an excited state and so is going to emit a photon, we cannot say *when* it will emit the photon. It has a certain amplitude to emit the photon at any time, and we can predict only a probability for emission; we cannot predict the future exactly. This has given rise to all kinds of nonsense and questions on the meaning of freedom of will, and of the idea that the world is uncertain.

Of course we must emphasize that classical physics is also indeterminate, in a sense. It is usually thought that this indeterminacy, that we cannot predict the future, is an important quantum-mechanical thing, and this is said to explain the behavior of the mind, feelings of free will, etc. But if the world *were* classical—if the laws of mechanics were classical—it is not quite obvious that the mind would not feel more or less the same. It is true classically that if we knew the position and the velocity of every particle in the world, or in a box of gas, we could predict exactly what would happen. And therefore the classical world is deterministic. Suppose, however, that we have a finite accuracy and do not know *exactly* where just one atom is, say to one part in a billion. Then as it goes along it hits another atom, and because we did not know the position better than to one part in a billion, we find an even larger error in the position after the collision. And that is amplified, of course, in the next collision, so that if we start with only a tiny error it rapidly magnifies to a very great uncertainty. To give an example: if water falls over a dam, it splashes. If we stand nearby, every now and then a drop will land on our nose. This appears to be completely random, yet such a behavior would be predicted by purely classical laws. The exact position of all the drops depends upon the precise wiggles of the water before it goes over the dam. How? The tiniest irregularities are magnified in falling, so that we get complete randomness. Ob-

viously, we cannot really predict the position of the drops unless we know the motion of the water *absolutely exactly*.

Speaking more precisely, given an arbitrary accuracy, no matter how precise, one can find a time long enough that we cannot make predictions valid for that long a time. Now the point is that this length of time is not very large. It is not that the time is millions of years if the accuracy is one part in a billion. The time goes, in fact, only logarithmically with the error, and it turns out that in only a very, very tiny time we lose all our information. If the accuracy is taken to be one part in billions and billions and billions—no matter how many billions we wish, provided we do stop somewhere—then we can find a time less than the time it took to state the accuracy—after which we can no longer predict what is going to happen! It is therefore not fair to say that from the apparent freedom and indeterminacy of the human mind, we should have realized that classical “deterministic” physics could not ever hope to understand it, and to welcome quantum mechanics as a release from a “completely mechanistic” universe. For already in classical mechanics there was indeterminability from a practical point of view.

The Kinetic Theory of Gases

39-1 Properties of matter

With this chapter we begin a new subject which will occupy us for some time. It is the first part of the analysis of the properties of matter from the physical point of view, in which, recognizing that matter is made out of a great many atoms, or elementary parts, which interact electrically and obey the laws of mechanics, we try to understand why various aggregates of atoms behave the way they do.

It is obvious that this is a difficult subject, and we emphasize at the beginning that it is in fact an *extremely* difficult subject, and that we have to deal with it differently than we have dealt with the other subjects so far. In the case of mechanics and in the case of light, we were able to begin with a precise statement of some laws, like Newton's laws, or the formula for the field produced by an accelerating charge, from which a whole host of phenomena could be essentially understood, and which would produce a basis for our understanding of mechanics and of light from that time on. That is, we may learn more later, but we do not learn different physics, we only learn better methods of mathematical analysis to deal with the situation.

We cannot use this approach effectively in studying the properties of matter. We can discuss matter only in a most elementary way; it is much too complicated a subject to analyze directly from its specific basic laws, which are none other than the laws of mechanics and electricity. But these are a bit too far away from the properties we wish to study; it takes too many steps to get from Newton's laws to the properties of matter, and these steps are, in themselves, fairly complicated. We will now start to take some of these steps, but while many of our analyses will be quite accurate, they will eventually get less and less accurate. We will have only a rough understanding of the properties of matter.

One of the reasons that we have to perform the analysis so imperfectly is that the mathematics of it requires a deep understanding of the theory of probability; we are not going to want to know where every atom is actually moving, but rather, how many move here and there on the average, and what the odds are for different effects. So this subject involves a knowledge of the theory of probability, and our mathematics is not yet quite ready and we do not want to strain it too hard.

Secondly, and more important from a physical standpoint, the actual behavior of the atoms is not according to classical mechanics, but according to quantum mechanics, and a correct understanding of the subject cannot be attained until we understand quantum mechanics. Here, unlike the case of billiard balls and automobiles, the difference between the classical mechanical laws and the quantum-mechanical laws is very important and very significant, so that many things that we will deduce by classical physics will be fundamentally incorrect. Therefore there will be certain things to be partially unlearned; however, we shall indicate in every case when a result is incorrect, so that we will know just where the "edges" are. One of the reasons for discussing quantum mechanics in the preceding chapters was to give an idea as to why, more or less, classical mechanics is incorrect in the various directions.

Why do we deal with the subject now at all? Why not wait a half a year, or a year, until we know the mathematics of probability better, and we learn a little quantum mechanics, and then we can do it in a more fundamental way? The answer is that it is a difficult subject, and the best way to learn is to do it slowly! The first thing to do is to get some idea, more or less, of what ought to happen in

39-1 Properties of matter

39-2 The pressure of a gas

39-3 Compressibility of radiation

39-4 Temperature and kinetic energy

39-5 The ideal gas law

different circumstances, and then, later, when we know the laws better, we will formulate them better.

Anyone who wants to analyze the properties of matter in a real problem might want to start by writing down the fundamental equations and then try to solve them mathematically. Although there are people who try to use such an approach, these people are the failures in this field; the real successes come to those who start from a *physical* point of view, people who have a rough idea where they are going and then begin by making the right kind of approximations, knowing what is big and what is small in a given complicated situation. These problems are so complicated that even an elementary understanding, although inaccurate and incomplete, is worth while having, and so the subject will be one that we shall go over again and again, each time with more and more accuracy, as we go through our course in physics.

Another reason for beginning the subject right now is that we have already used many of these ideas in, for example, chemistry, and we have even heard of some of them in high school. It is interesting to know the physical basis for these things.

As an interesting example, we all know that equal volumes of gases, at the same pressure and temperature, contain the same number of molecules. The law of multiple proportions, that when two gases combine in a chemical reaction the volumes needed always stand in simple integral proportions, was understood ultimately by Avogadro to mean that equal volumes have equal numbers of atoms. Now *why* do they have equal numbers of atoms? Can we deduce from Newton's laws that the number of atoms should be equal? We shall address ourselves to that specific matter in this chapter. In succeeding chapters, we shall discuss various other phenomena involving pressures, volumes, temperature, and heat.

We shall also find that the subject can be attacked from a nonatomic point of view, and that there are many interrelationships of the properties of substances. For instance, when we compress something, it heats; if we heat it, it expands. There is a relationship between these two facts which can be deduced independently of the machinery underneath. This subject is called *thermodynamics*. The deepest understanding of thermodynamics comes, of course, from understanding the actual machinery underneath, and that is what we shall do: we shall take the atomic viewpoint from the beginning and use it to understand the various properties of matter and the laws of thermodynamics.

Let us, then, discuss the properties of gases from the standpoint of Newton's laws of mechanics.

39-2 The pressure of a gas

First, we know that a gas exerts a pressure, and we must clearly understand what this is due to. If our ears were a few times more sensitive, we would hear a perpetual rushing noise. Evolution has not developed the ear to that point, because it would be useless if it were so much more sensitive—we would hear a perpetual racket. The reason is that the eardrum is in contact with the air, and air is a lot of molecules in perpetual motion and these bang against the eardrums. In banging against the eardrums they make an irregular tattoo—boom, boom, boom—which we do not hear because the atoms are so small, and the sensitivity of the ear is not quite enough to notice it. The result of this perpetual bombardment is to push the drum away, but of course there is an equal perpetual bombardment of atoms on the other side of the eardrum, so the net force on it is zero. If we were to take the air away from one side, or change the relative amounts of air on the two sides, the eardrum would then be pushed one way or the other, because the amount of bombardment on one side would be greater than on the other. We sometimes feel this uncomfortable effect when we go up too fast in an elevator or an airplane, especially if we also have a bad cold (when we have a cold, inflammation closes the tube which connects the air on the inside of the eardrum with the outside air through the throat, so that the two pressures cannot readily equalize).

In considering how to analyze the situation quantitatively, we imagine that we have a volume of gas in a box, at one end of which is a piston which can be moved (Fig. 39-1). We would like to find out what force on the piston results from the fact that there are atoms in this box. The volume of the box is V , and as the atoms move around inside the box with various velocities they bang against the piston. Suppose there is nothing, a vacuum, on the outside of the piston. What of it? If the piston were left alone, and nobody held onto it, each time it got banged it would pick up a little momentum and it would gradually get pushed out of the box. So in order to keep it from being pushed out of the box, we have to hold it with a force F . The problem is, how much force? One way of expressing the force is to talk about the force per unit area: if A is the area of the piston, then the force on the piston will be written as a number times the area. We define the pressure, then, as equal to the force that we have to apply on a piston, divided by the area of the piston:

$$P = F/A. \quad (39.1)$$

To make sure we understand the idea (we have to derive it for another purpose anyway), the differential work dW done on the gas in compressing it by moving the piston in a differential amount $-dx$ would be the force times the distance that we compress it, which, according to (39.1), would be the pressure times the area, times the distance, which is equal to minus the pressure times the change in the volume:

$$dW = F(-dx) = -PA dx = -P dV. \quad (39.2)$$

(The area A times the distance dx is the volume change.) The minus sign is there because, as we compress it, we *decrease* the volume; if we think about it we can see that if a gas is compressed, work is done *on* it.

How much force do we have to apply to balance the banging of the molecules? The piston receives from each collision a certain amount of momentum. A certain amount of momentum per second will pour into the piston, and it will start to move. To keep it from moving, we must pour back into it the same amount of momentum per second from our force. Of course, the force *is* the amount of momentum per second that we must pour in. There is another way to put it: if we let go of the piston it will pick up speed because of the bombardments; with each collision we get a little more speed, and the speed thus accelerates. The rate at which the piston picks up speed, or accelerates, is proportional to the force on it. So we see that the force, which we already have said is the pressure times the area, is equal to the momentum per second delivered to the piston by the colliding molecules.

To calculate the momentum per second is easy—we can do it in two parts: first, we find the momentum delivered to the piston by one particular atom in a collision with the piston, then we have to multiply by the number of collisions per second that the atoms have with the wall. The force will be the product of these two factors. Now let us see what the two factors are: In the first place, we shall suppose that the piston is a perfect “reflector” for the atoms. If it is not, the whole theory is wrong, and the piston will start to heat up and things will change, but eventually, when equilibrium has set in, the net result is that the collisions are effectively perfectly elastic. On the average, every particle that comes in leaves with the same energy. So we shall imagine that the gas is in a steady condition, and we lose no energy to the piston because the piston is standing still. In those circumstances, if a particle comes in with a certain speed, it comes out with the same speed and, we will say, with the same mass.

If \mathbf{v} is the velocity of an atom, and v_x is the x -component of \mathbf{v} , then mv_x is the x -component of momentum “in”; but we also have an equal component of momentum “out,” and so the total momentum delivered to the piston by the particle, in one collision, is $2mv_x$, because it is “reflected.”

Now, we need the number of collisions made by the atoms in a second, or in a certain amount of time dt ; then we divide by dt . How many atoms are hitting? Let us suppose that there are N atoms in the volume V , or $n = N/V$ in each unit volume. To find how many atoms hit the piston, we note that, given a certain

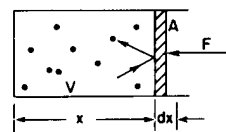


Fig. 39-1. Atoms of a gas in a box with a frictionless piston.

amount of time t , if a particle has a certain velocity toward the piston it will hit during the time t , provided it is close enough. If it is too far away, it goes only part way toward the piston in the time t , but does not reach the piston. Therefore it is clear that only those molecules which are within a distance $v_x t$ from the piston are going to hit the piston in the time t . Thus the number of collisions in a time t is equal to the number of atoms which are in the region within a distance $v_x t$, and since the area of the piston is A , the volume occupied by the atoms which are going to hit the piston is $v_x t A$. But the number of atoms that are going to hit the piston is that volume times the number of atoms per unit volume, $n v_x t A$. Of course we do not want the number that hit in a time t , we want the number that hit per second, so we divide by the time t , to get $n v_x A$. (This time t could be made very short; if we feel we want to be more elegant, we call it dt , then differentiate, but it is the same thing.)

So we find that the force is

$$F = n v_x A \cdot 2 m v_x. \quad (39.3)$$

See, the force is proportional to the area, if we keep the particle density fixed as we change the area! The pressure is then

$$P = 2 n m v_x^2. \quad (39.4)$$

Now we notice a little trouble with this analysis: First, all the molecules do not have the same velocity, and they do not move in the same direction. So, all the v_x^2 's are different! So what we must do, of course, is to take an *average* of the v_x^2 's, since each one makes its own contribution. What we want is the square of v_x , averaged over all the molecules:

$$P = n m \langle v_x^2 \rangle. \quad (39.5)$$

Did we forget to include the factor 2? No; of all the atoms, only half are headed toward the piston. The other half are headed the other way, and if we take $\langle v_x^2 \rangle$, we are averaging the *negative* v_x 's squared, as well as the positive v_x 's. So when we just take $\langle v_x^2 \rangle$, without looking, we are getting twice as much as we want. The average of v_x^2 , for positive v_x , is equal to the average of v_x^2 for *all* v_x , times one half.

Now as the atoms bounce around, it is clear that there is nothing special about the "x-direction"; the atoms may also be moving up and down, back and forth, in and out. Therefore it is going to be true that $\langle v_x^2 \rangle$, the average motion of the atoms in one direction, and the average in the other two directions, are all going to be equal:

$$\langle v_x^2 \rangle = \langle v_y^2 \rangle = \langle v_z^2 \rangle. \quad (39.6)$$

It is only a matter of rather tricky mathematics to notice, therefore, that they are each equal to one-third of their sum, which is of course the square of the magnitude of the velocity:

$$\langle v_x^2 \rangle = \frac{1}{3} \langle v_x^2 + v_y^2 + v_z^2 \rangle = \langle v^2 \rangle / 3. \quad (39.7)$$

This has the advantage that we do not have to worry about any particular direction, and so we write our pressure formula again in this form:

$$P = \left(\frac{2}{3}\right) n \langle m v^2 / 2 \rangle. \quad (39.8)$$

The reason we wrote the last factor as $\langle m v^2 / 2 \rangle$ is that this is the *kinetic energy* of the center-of-mass motion of the molecule. We find, therefore, that

$$P V = N \left(\frac{2}{3}\right) \langle m v^2 / 2 \rangle. \quad (39.9)$$

With this equation we can calculate how much the pressure is, if we know the speeds.

As a very simple example let us take helium gas, or any other gas, like mercury vapor, or potassium vapor of high enough temperature, or argon, in which all the molecules are single atoms, for which we may suppose that there is no internal

motion in the atom. If we had a complex molecule, there might be some internal motion, mutual vibrations, or something. We suppose that we may disregard that; this is actually a serious matter that we will have to come back to, but it turns out to be all right. We suppose that the internal motion of the atoms can be disregarded, and therefore, for this purpose, that the kinetic energy of the center-of-mass motion is all the energy there is. So for a monatomic gas, the kinetic energy is the total energy. In general, we are going to call U the total energy (it is sometimes called the total *internal* energy—we may wonder why, since there is no *external* energy to a gas), i.e., all the energy of all the molecules in the gas, or the object, whatever it is.

For a monatomic gas we will suppose that the total energy U is equal to a number of atoms times the average kinetic energy of each, because we are disregarding any possibility of excitation or motion inside the atoms themselves. Then, in these circumstances, we would have

$$PV = \frac{2}{3}U. \quad (39.10)$$

Incidentally, we can stop here and find the answer to the following question: Suppose that we take a can of gas and compress the gas slowly, how much pressure do we need to squeeze the volume down? It is easy to find out, since the pressure is $\frac{2}{3}$ the energy divided by V . As we squeeze it down, we do work on the gas and we thereby increase the energy U . So we are going to have some kind of a differential equation: If we start out in a given circumstance with a certain energy and a certain volume, we then know the pressure. Now we start to squeeze, but the moment we do, the energy U increases and the volume V decreases, so the pressure goes up.

So, we have to solve a differential equation, and we will solve it in a moment. We must first emphasize, however, that as we are compressing this gas, we are supposing that all the work goes into increasing the energy of the atoms inside. We may ask, "Isn't that necessary? Where else could it go?" It turns out that it *can* go another place. There are what we call "heat leaks" through the walls: the hot (i.e., fast-moving) atoms that bombard the walls, heat the walls, and energy goes away. We shall suppose for the present that this is not the case.

For somewhat wider generality, although we are still making some very special assumptions about our gas, we shall write, not $PV = \frac{2}{3}U$, but

$$PV = (\gamma - 1)U. \quad (39.11)$$

It is written $(\gamma - 1)$ times U for conventional reasons, because we will deal with a few other cases later where the number in front of U will not be $\frac{2}{3}$, but will be a different number. So, in order to do the thing in general, we call it $\gamma - 1$, because people have been calling it that for almost one hundred years. This γ , then, is $\frac{5}{3}$, because $\frac{5}{3} - 1$ is $\frac{2}{3}$ for a monatomic gas like helium.

We have already noticed that when we compress a gas the work done is $-P dV$. A compression in which there is no heat energy added or removed is called an *adiabatic* compression, from the Greek *a* (not) + *dia* (through) + *bainein* (to go). (The word adiabatic is used in physics in several ways, and it is sometimes hard to see what is common about them.) That is, for an adiabatic compression all the work done goes into changing the internal energy. That is the key—that there are no other losses of energy—for then we have $P dV = -dU$. But since $U = PV/(\gamma - 1)$, we may write

$$dU = (P dV + V dP)/(\gamma - 1). \quad (39.12)$$

So we have $P dV = -(P dV + V dP)/(\gamma - 1)$, or, rearranging the terms, $\gamma P dV = -V dP$, or

$$(\gamma dV/V) + (dP/P) = 0. \quad (39.13)$$

Fortunately, assuming that γ is constant, as it is for a monatomic gas, we can integrate this: it gives $\gamma \ln V + \ln P = \ln C$, where $\ln C$ is the constant of integration. If we take the exponential of both sides, we get the law

$$PV^\gamma = C \text{ (a constant)}. \quad (39.14)$$

In other words, under adiabatic conditions, where the temperature rises as we compress because no heat is being lost, the pressure times the volume to the $\frac{5}{3}$ power is a constant for a monatomic gas! Although we derived it theoretically, this *is*, in fact, the way monatomic gases behave experimentally.

39-3 Compressibility of radiation

We may give one other example of the kinetic theory of a gas, one which is not used in chemistry so much, but is used in astronomy. We have a large number of photons in a box in which the temperature is very high. (The box is, of course, the gas in a very hot star. The sun is not hot enough; there are still too many atoms, but at still higher temperatures in certain very hot stars, we may neglect the atoms and suppose that the only objects that we have in the box are photons.) Now then, a photon has a certain momentum \mathbf{p} . (We always find that we are in terrible trouble when we do kinetic theory: p is the pressure, but p is the momentum; v is the volume, but v is the velocity; T is the temperature, but T is the kinetic energy or the time or the torque; one must keep one's wits about one!) This \mathbf{p} is momentum, it is a vector. Going through the same analysis as before, it is the x -component of the vector \mathbf{p} which generates the "kick," and twice the x -component of the vector \mathbf{p} is the momentum which is given in the kick. Thus $2p_x$ replaces $2mv_x$, and in evaluating the number of collisions, v_x is still v_x , so when we get all the way through, we find that the pressure in Eq. (39.4) is, instead,

$$P = 2np_xv_x. \quad (39.15)$$

Then, in the averaging, it becomes n times the average of p_xv_x (the same factor of 2) and, finally, putting in the other two directions, we find

$$PV = N\langle \mathbf{p} \cdot \mathbf{v} \rangle / 3. \quad (39.16)$$

This checks with the formula (39.9), because the momentum is $m\mathbf{v}$; it is a little more general, that is all. The pressure times the volume is the total number of atoms times $\frac{1}{3}\langle \mathbf{p} \cdot \mathbf{v} \rangle$, averaged.

Now, for photons, what is $\mathbf{p} \cdot \mathbf{v}$? The momentum and the velocity are in the same direction, and the velocity is the speed of light, so this is the momentum of each of the objects, times the speed of light. The momentum times the speed of light of every photon is its energy: $E = pc$, so these terms are the *energies* of each of the photons, and we should, of course, take an average energy, times the number of photons. So we have $\frac{1}{3}$ of the energy inside the gas:

$$PV = U/3 \text{ (photon gas)}. \quad (39.17)$$

For photons, then, since we have $\frac{1}{3}$ in front, $(\gamma - 1)$ in (39.11) is $\frac{1}{3}$, or $\gamma = \frac{4}{3}$, and we have discovered that radiation in a box obeys the law

$$PV^{4/3} = C. \quad (39.18)$$

So we know the compressibility of radiation! That is what is used in an analysis of the contribution of radiation pressure in a star, that is how we calculate it, and how it changes when we compress it. What wonderful things are already within our power!

39-4 Temperature and kinetic energy

So far we have not dealt with *temperature*; we have purposely been avoiding the temperature. As we compress a gas, we know that the energy of the molecules increases, and we are used to saying that the gas gets hotter; we would like to understand what this has to do with the temperature. If we try to do the experiment, not adiabatically but at what we call *constant temperature*, what are we doing? We know that if we take two boxes of gas and let them sit next to each other long enough, even if at the start they were at what we call different temperatures, they

will in the end come to the same temperature. Now what does that mean? That means that they get to a condition that they would get to if we left them alone long enough! What we mean by equal temperature is just that—the final condition when things have been sitting around interacting with each other long enough.

Let us consider, now, what happens if we have two gases in containers separated by a movable piston as in Fig. 39-2 (just for simplicity we shall take two monatomic gases, say helium and neon). In container (1) the atoms have mass m_1 , velocity v_1 , and there are n_1 per unit volume, and in the other container the atoms have mass m_2 , velocity v_2 , there are n_2 atoms per unit volume. What are the conditions for equilibrium?

Obviously, the bombardment from the left side must be such that it moves the piston to the right and compresses the other gas until its pressure builds up, and the thing will thus slosh back and forth, and will gradually come to rest at a place where the pressures are equal on both sides. So we can arrange that the pressures are equal; that just means that the internal energies per unit volume are equal, or that the numbers n times the average kinetic energies on each side are equal. What we have to try to prove, eventually, is that the *numbers themselves* are equal. So far, all we know is that the numbers times the kinetic energies are equal,

$$n_1 \langle m_1 v_1^2 / 2 \rangle = n_2 \langle m_2 v_2^2 / 2 \rangle,$$

from (39.8), because the pressures are equal. We must realize that this is not the only condition over the long run, but something else must happen more slowly as the true complete equilibrium corresponding to equal temperatures sets in.

To see the idea, suppose that the pressure on the left side were developed by having a very high density but a low velocity. By having a large n and a small v , we can get the same pressure as by having a small n and a large v . The atoms may be moving slowly but be packed nearly solidly, or there may be fewer but they are hitting harder. Will it stay like that forever? At first we might think so, but then we think again and find we have forgotten one important point. That is, that the intermediate piston does not receive a steady pressure; it wiggles, just like the eardrum that we were first talking about, because the bangings are not absolutely uniform. There is not a perpetual, steady pressure, but a tattoo—the pressure varies, and so the thing jiggles. Suppose that the atoms on the right side are not jiggling much, but those on the left are few and far between and very energetic. The piston will, now and then, get a big impulse from the left, and will be driven against the slow atoms on the right, giving them more speed. (As each atom collides with the piston, it either gains or loses energy, depending upon whether the piston is moving one way or the other when the atom strikes it.) So, as a result of the collisions, the piston finds itself jiggling, jiggling, jiggling, and this shakes the other gas—it gives energy to the other atoms, and they build up faster motions, until they balance the jiggling that the piston is giving to them. The system comes to some equilibrium where the piston is moving at such a mean square speed that it picks up energy from the atoms at about the same rate as it puts energy back into them. So the piston picks up a certain mean irregularity in speed, and it is our problem to find it. When we do find it, we can solve our problem better, because the gases will adjust their velocities until the rate at which they are trying to pour energy into each other through the piston will become equal.

It is quite difficult to figure out the details of the piston in this particular circumstance; although it is ideally simple to understand, it turns out to be a little harder to analyze. Before we analyze that, let us analyze another problem in which we have a box of gas but now we have two different kinds of molecules in it, having masses m_1 and m_2 , velocities v_1 and v_2 , and so forth; there is now a much more intimate relationship. If all of the No. 2 molecules are standing still, that condition is not going to last, because they get kicked by the No. 1 molecules and so pick up speed. If they are all going much faster than the No. 1 molecules, then maybe that will not last either—they will pass the energy back to the No. 1 molecules. So when both gases are in the same box, the problem is to find the rule that determines the relative speeds of the two.

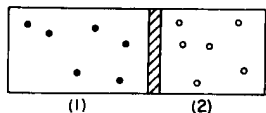


Fig. 39-2. Atoms of two different monatomic gases are separated by a movable piston.

This is still a very difficult problem, but we will solve it as follows. First we consider the following sub-problem (again this is one of those cases where—never mind the derivation—in the end the result is very simple to remember, but the derivation is just ingenious). Let us suppose that we have two molecules, of different mass, colliding, and that the collision is viewed on the center-of-mass (CM) system. In order to remove a complication, we look at the collision in the CM. As we know from the laws of collision, by the conservation of momentum and energy, after the molecules collide the only way they can move is such that each maintains its own original speed—and they just change their *direction*. So we have an average collision that looks like that in Fig. 39-3. Suppose, for a moment, that we watch all the collisions with the CM at rest. Suppose we imagine that they are all initially moving horizontally. Of course, after the first collision some of them are moving at an angle. In other words, if they were all going horizontally, then at least some would later be moving vertically. Now in some other collision, they would be coming in from another direction, and then they would be twisted at still another angle. So even if they were completely organized in the beginning, they would get sprayed around at all angles, and then the sprayed ones would get sprayed some more, and sprayed some more, and sprayed some more. Ultimately, what will be the distribution? *Answer: It will be equally likely to find any pair moving in any direction in space.* After that further collisions could not change the distribution.

They are equally likely to go in all directions, but how do we say that? There is of course *no* likelihood that they will go in any specific direction, because a specific direction is too exact, so we have to talk about per unit “something.” The idea is that any area on a sphere centered at a collision point will have just as many molecules going through it as go through any other equal area on the sphere. So the result of the collisions will be to distribute the directions so that equal areas on a sphere will have equal probabilities.

Incidentally, if we just want to discuss the original direction and some other direction an angle θ from it, it is an interesting property that the differential area of a sphere of unit radius is $\sin \theta d\theta$ times 2π , and that is the same as the differential of $\cos \theta$. So what it means is that the cosine of the angle θ between any two directions is equally likely to be anything from -1 to $+1$.

Next, we have to worry about the actual case, where we do not have the collision in the CM system, but we have two atoms which are coming together with vector velocities \mathbf{v}_1 and \mathbf{v}_2 . What happens now? We can analyze this collision with the vector velocities \mathbf{v}_1 and \mathbf{v}_2 in the following way: We first say that there is a certain CM; the velocity of the CM is given by the “average” velocity, with weights proportional to the masses, so the velocity of the CM is $\mathbf{v}_{CM} = (m_1\mathbf{v}_1 + m_2\mathbf{v}_2)/(m_1 + m_2)$. If we watch this collision in the CM system, then we see a collision just like that in Fig. 39-3, with a certain relative velocity \mathbf{w} coming in. The relative velocity is just $\mathbf{v}_1 - \mathbf{v}_2$. Now the idea is that, first, the whole CM is moving, and in the CM there is a relative velocity \mathbf{w} , and the molecules collide and come off in some new direction. All this happens while the CM keeps right on moving, without any change.

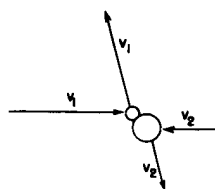


Fig. 39-3. A collision between unequal molecules, viewed in the CM system.

Now then, what is the distribution resulting from this? From our previous argument we conclude this: that at equilibrium, *all directions for \mathbf{w} are equally likely, relative to the direction of the motion of the CM.** There will be no particular correlation, in the end, between the direction of the motion of the relative velocity and that of the motion of the CM. Of course, if there were, the collisions would spray it about, so it is all sprayed around. So the cosine of the angle between \mathbf{w} and \mathbf{v}_{CM} is zero on the average. That is,

$$\langle \mathbf{w} \cdot \mathbf{v}_{CM} \rangle = 0. \quad (39.19)$$

* This argument, which was the one used by Maxwell, involves some subtleties. Although the conclusion is correct, the result does *not* follow purely from the considerations of symmetry that we used before, since, by going to a reference frame moving through the gas, we may find a distorted velocity distribution. We have not found a simple proof of this result.

But $\mathbf{w} \cdot \mathbf{v}_{CM}$ can be expressed in terms of \mathbf{v}_1 and \mathbf{v}_2 as well:

$$\begin{aligned}\mathbf{w} \cdot \mathbf{v}_{CM} &= \frac{(\mathbf{v}_1 - \mathbf{v}_2) \cdot (m_1 \mathbf{v}_1 + m_2 \mathbf{v}_2)}{m_1 + m_2} \\ &= \frac{(m_1 v_1^2 - m_2 v_2^2) + (m_2 - m_1)(\mathbf{v}_1 \cdot \mathbf{v}_2)}{m_1 + m_2}.\end{aligned}\quad (39.20)$$

First, let us look at the $\mathbf{v}_1 \cdot \mathbf{v}_2$; what is the average of $\mathbf{v}_1 \cdot \mathbf{v}_2$? That is, what is the average of the component of velocity of one molecule in the direction of another? Surely there is just as much likelihood of finding any given molecule moving one way as another. *The average of the velocity \mathbf{v}_2 in any direction is zero.* Certainly, then, in the direction of \mathbf{v}_1 , \mathbf{v}_2 has zero average. So, the average of $\mathbf{v}_1 \cdot \mathbf{v}_2$ is zero! Therefore, we conclude that the average of $m_1 v_1^2$ must be equal to the average of $m_2 v_2^2$. That is, *the average kinetic energy of the two must be equal:*

$$\frac{1}{2} m_1 v_1^2 = \frac{1}{2} m_2 v_2^2. \quad (39.21)$$

If we have two kinds of atoms in a gas, it can be shown, and we presume to have shown it, that the average of the kinetic energy of one is the same as the average of the kinetic energy of the other, when they are both in the same gas in the same box in equilibrium. That means that the heavy ones will move slower than the light ones; this is easily shown by experimentation with “atoms” of different masses in an air trough.

Now we would like to go one step further, and say that if we have two different gases *separated* in a box, they will also have equal average kinetic energy when they have finally come to equilibrium, even though they are not in the same box. We can make the argument in a number of ways. One way is to argue that if we have a fixed partition with a tiny hole in it (Fig. 39-4) so that one gas could leak out through the holes while the other could not, because the molecules are too big, and these had attained equilibrium, then we know that in one part, where they are mixed, they have the same average kinetic energy, but some come through the hole without loss of kinetic energy, so the average kinetic energy in the pure gas and in the mixture must be the same. That is not too satisfactory, because maybe there are no holes, for this kind of molecule, that separate one kind from the other.

Let us now go back to the piston problem. We can give an argument which shows that the kinetic energy of this piston must also be $\frac{1}{2} m_2 v_2^2$. Actually, that would be the kinetic energy due to the purely horizontal motion of the piston, so, forgetting its up and down motion, it will have to be the same as $\frac{1}{2} m_2 v_{2x}^2$. Likewise, from the equilibrium on the other side, we can prove that the kinetic energy of the piston is $\frac{1}{2} m_1 v_{1x}^2$. Although this is not in the middle of the gas, but is on one side of the gas, we can still make the argument, although it is a little more difficult, that the average kinetic energy of the piston and of the gas molecules are equal as a result of all the collisions.

If this still does not satisfy us, we may make an artificial example by which the equilibrium is generated by an object which can be hit on all sides. Suppose that we have a short rod with a ball on each end sticking through the piston, on a frictionless sliding universal joint. Each ball is round, like one of the molecules, and can be hit on all sides. This whole object has a certain total mass, m . Now, we have the gas molecules with mass m_1 and mass m_2 as before. The result of the collisions, by the analysis that was made before, is that the kinetic energy of m because of collisions with the molecules on one side must be $\frac{1}{2} m_1 v_1^2$, on the average. Likewise, because of the collisions with molecules on the other side, it has to be $\frac{1}{2} m_2 v_2^2$ on the average. So, therefore, both sides have to have the *same* kinetic energy when they are in thermal equilibrium. So, although we only proved it for a mixture of gases, it is easily extended to the case where there are two different, separate gases at the same temperature.

Thus *when we have two gases at the same temperature, the mean kinetic energy of the CM motions are equal.*

The mean molecular kinetic energy is a property only of the “temperature.” Being a property of the “temperature,” *and not of the gas*, we can use it as a *definition* of the temperature. The mean kinetic energy of a molecule is thus some

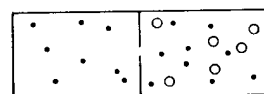


Fig. 39-4. Two gases in a box with a semipermeable membrane.

function of the temperature. But who is to tell us what scale to use for the temperature? We may arbitrarily *define* the scale of temperature so that the mean energy is linearly proportional to the temperature. The best way to do it would be to call the mean energy itself “the temperature.” That would be the simplest possible function. Unfortunately, the scale of temperature has been chosen differently, so instead of calling it temperature directly we use a constant conversion factor between the energy of a molecule and a degree of absolute temperature called a degree Kelvin. The constant of proportionality is $k = 1.38 \times 10^{-23}$ joule for every degree Kelvin.* So if T is absolute temperature, our definition says that the mean molecular kinetic energy is $\frac{3}{2} kT$. (The $\frac{3}{2}$ is put in as a matter of convenience, so as to get rid of it somewhere else.)

We point out that the kinetic energy associated with the component of motion in any particular direction is only $\frac{1}{2} kT$. The three independent directions that are involved make it $\frac{3}{2} kT$.

39-5 The ideal gas law

Now, of course, we can put our definition of temperature into Eq. (39.9) and so find the law for the pressure of gases as a function of the temperature: it is that the pressure times the volume is equal to the total number of atoms times the universal constant k , times the temperature:

$$PV = NkT. \quad (39.22)$$

Furthermore, at the same temperature and pressure and volume, the *number of atoms* is determined; it too is a universal constant! So equal volumes of different gases, at the same pressure and temperature, have the same number of molecules, because of Newton’s laws. That is an amazing conclusion!

In practice, when dealing with molecules, because the numbers are so large, the chemists have artificially chosen a specific number, a very large number, and called it something else. They have a number which they call a *mole*. A mole is merely a handy number. Why they did not choose 10^{24} objects, so it would come out even, is a historical question. They happened to choose, for the convenient number of objects on which they standardize, $N_0 = 6.02 \times 10^{23}$ objects, and this is called a mole of objects. So instead of measuring the number of molecules in units, they measure in terms of numbers of moles.† In terms of N_0 we can write the number of moles, times the number of atoms in a mole, times kT , and if we want to, we can take the number of atoms in a mole times k , which is a mole’s worth of k , and call it something else, and we do—we call it R . A mole’s worth of k is 8.317 joules: $R = N_0 k = 8.317 \text{ j} \cdot \text{mole}^{-1} \cdot ^\circ\text{K}^{-1}$. Thus we also find the gas law written as the number of moles (also called N) times RT , or the number of atoms, times kT :

$$PV = NRT. \quad (39.23)$$

It is the same thing, just a different scale for measuring numbers. We use 1 as a unit, and chemists use 6×10^{23} as a unit!

We now make one more remark about our gas law, and that has to do with the law for objects other than monatomic molecules. We have dealt only with the CM motion of the atoms of a monatomic gas. What happens if there are forces present? First, consider the case that the piston is held by a horizontal spring, and there are forces on it. The exchange of jiggling motion between atoms and piston at any moment does not depend on where the piston is at that moment, of course. The equilibrium conditions are the same. No matter where the piston is, its speed of motion must be such that it passes energy to the molecules in just

* The centigrade scale is just this Kelvin scale with a zero chosen at 273.16 °K, so $T = 273.16 + \text{centigrade temperature}$.

† What the chemists call molecular weights are the masses in grams of a mole of a molecule. The mole is defined so that the mass of a mole of carbon atoms of isotope 12 (i.e., having 6 protons and 6 neutrons in the nucleus) is exactly 12 grams.

the right way. So it makes no difference about the spring. The *speed* at which the piston has to move, on the average, is the same. So our theorem, that the mean value of the kinetic energy in one direction is $\frac{1}{2}kT$, is true whether there are forces present or not.

Consider, for example, a diatomic molecule composed of atoms m_A and m_B . What we have proved is that the motion of the CM of part A and that of part B are such that $\langle \frac{1}{2}m_A v_A^2 \rangle = \langle \frac{1}{2}m_B v_B^2 \rangle = \frac{3}{2}kT$. How can this be, if they are held together? Although they are held together, when they are spinning and turning in there, when something hits them, exchanging energy with them, *the only thing that counts is how fast they are moving*. That alone determines how fast they exchange energy in collisions. At the particular instant, the force is not an essential point. Therefore the same principle is right, even when there are forces.

Let us prove, finally, that the gas law is consistent also with a disregard of the internal motion. We did not really include the internal motions before; we just treated a monatomic gas. But we shall now show that an entire object, considered as a single body of total mass M , has a velocity of the CM such that

$$\frac{1}{2}Mv_{CM}^2 = \frac{3}{2}kT. \quad (39.24)$$

In other words, we can consider either the separate pieces or the whole thing! Let us see the reason for that: The mass of the diatomic molecule is $M = m_A + m_B$, and the velocity of the center of mass is equal to $\mathbf{v}_{CM} = (m_A \mathbf{v}_A + m_B \mathbf{v}_B)/M$. Now we need $\langle v_{CM}^2 \rangle$. If we square \mathbf{v}_{CM} , we get

$$v_{CM}^2 = \frac{m_A^2 v_A^2 + 2m_A m_B \mathbf{v}_A \cdot \mathbf{v}_B + m_B^2 v_B^2}{M^2}.$$

Now we multiply $\frac{1}{2}M$ and take the average, and thus we get

$$\begin{aligned} \frac{1}{2}Mv_{CM}^2 &= \frac{m_A \frac{3}{2}kT + 2m_A m_B \langle \mathbf{v}_A \cdot \mathbf{v}_B \rangle + m_B \frac{3}{2}kT}{M} \\ &= \frac{3}{2}kT + \frac{2m_A m_B \langle \mathbf{v}_A \cdot \mathbf{v}_B \rangle}{M}. \end{aligned}$$

(We have used the fact that $(m_A + m_B)/M = 1$.) Now what is $\langle \mathbf{v}_A \cdot \mathbf{v}_B \rangle$? (It had better be zero!) To find out, let us use our assumption that the relative velocity, $\mathbf{w} = \mathbf{v}_A - \mathbf{v}_B$ is not any more likely to point in one direction than in another—that is, that its average component in any direction is zero. Thus we assume that

$$\langle \mathbf{w} \cdot \mathbf{v}_{CM} \rangle = 0.$$

But what is $\mathbf{w} \cdot \mathbf{v}_{CM}$? It is

$$\begin{aligned} \mathbf{w} \cdot \mathbf{v}_{CM} &= \frac{(\mathbf{v}_A - \mathbf{v}_B) \cdot (m_A \mathbf{v}_A + m_B \mathbf{v}_B)}{M} \\ &= \frac{m_A v_A^2 + (m_B - m_A)(\mathbf{v}_A \cdot \mathbf{v}_B) - m_B v_B^2}{M}. \end{aligned}$$

Therefore, since $\langle m_A v_A^2 \rangle = \langle m_B v_B^2 \rangle$, the first and last terms cancel out on the average, and we are left with

$$(m_B - m_A) \langle \mathbf{v}_A \cdot \mathbf{v}_B \rangle = 0.$$

Thus if $m_A \neq m_B$, we find that $\langle \mathbf{v}_A \cdot \mathbf{v}_B \rangle = 0$, and therefore that the bodily motion of the entire molecule, regarded as a single particle of mass M , has a kinetic energy, on the average, equal to $\frac{3}{2}kT$.

Incidentally, we have also proved at the same time that the average kinetic energy of the *internal* motions of the diatomic molecule, disregarding the bodily motion of the CM, is $\frac{3}{2}kT$! For, the total kinetic energy of the parts of the molecule is $\frac{1}{2}m_A v_A^2 + \frac{1}{2}m_B v_B^2$, whose average is $\frac{3}{2}kT + \frac{3}{2}kT$, or $3kT$. The kinetic energy of the center-of-mass motion is $\frac{3}{2}kT$, so the average kinetic energy of the rotational and vibratory motions of the two atoms inside the molecule is the difference, $\frac{3}{2}kT$.

The theorem concerning the average energy of the CM motion is general: for any object considered as a whole, with forces present or no, for every independent direction of motion that there is, the average kinetic energy in that motion is $\frac{1}{2}kT$. These “independent directions of motion” are sometimes called the *degrees of freedom* of the system. The number of degrees of freedom of a molecule composed of r atoms is $3r$, since each atom needs three coordinates to define its position. The entire kinetic energy of the molecule can be expressed either as the sum of the kinetic energies of the separate atoms, or as the sum of the kinetic energy of the CM motion plus the kinetic energy of the internal motions. The latter can sometimes be expressed as a sum of rotational kinetic energy of the molecule and vibrational energy, but this is an approximation. Our theorem, applied to the r -atom molecule, says that the molecule will have, on the average, $3rkT/2$ joules of kinetic energy, of which $\frac{3}{2}kT$ is kinetic energy of the center-of-mass motion of the entire molecule, and the rest, $\frac{3}{2}(r - 1)kT$, is internal vibrational and rotational kinetic energy.

The Principles of Statistical Mechanics

40-1 The exponential atmosphere

We have discussed some of the properties of large numbers of intercolliding atoms. The subject is called kinetic theory, a description of matter from the point of view of collisions between the atoms. Fundamentally, we assert that the gross properties of matter should be explainable in terms of the motion of its parts.

We limit ourselves for the present to conditions of thermal equilibrium, that is, to a subclass of all the phenomena of nature. The laws of mechanics which apply just to thermal equilibrium are called *statistical mechanics*, and in this section we want to become acquainted with some of the central theorems of this subject.

We already have one of the theorems of statistical mechanics, namely, the mean value of the kinetic energy for any motion at the absolute temperature T is $\frac{1}{2}kT$ for each independent motion, i.e., for each degree of freedom. That tells us something about the mean square velocities of the atoms. Our objective now is to learn more about the positions of the atoms, to discover how many of them are going to be in different places at thermal equilibrium, and also to go into a little more detail on the distribution of the velocities. Although we have the mean square velocity, we do not know how to answer a question such as how many of them are going three times faster than the root mean square, or how many of them are going one-quarter of the root mean square speed. Or have they all the same speed exactly?

So, these are the two questions that we shall try to answer: How are the molecules distributed in space when there are forces acting on them, and how are they distributed in velocity?

It turns out that the two questions are completely independent, and that the distribution of velocities is always the same. We already received a hint of the latter fact when we found that the average kinetic energy is the same, $\frac{1}{2}kT$ per degree of freedom, no matter what forces are acting on the molecules. The distribution of the velocities of the molecules is independent of the forces, because the collision rates do not depend upon the forces.

Let us begin with an example: the distribution of the molecules in an atmosphere like our own, but without the winds and other kinds of disturbance. Suppose that we have a column of gas extending to a great height, and at thermal equilibrium—unlike our atmosphere, which as we know gets colder as we go up. We could remark that if the temperature differed at different heights, we could demonstrate lack of equilibrium by connecting a rod to some balls at the bottom (Fig. 40-1), where they would pick up $\frac{1}{2}kT$ from the molecules there and would shake, via the rod, the balls at the top and those would shake the molecules at the top. So, ultimately, of course, the temperature becomes the same at all heights in a gravitational field.

If the temperature is the same at all heights, the problem is to discover by what law the atmosphere becomes tenuous as we go up. If N is the total number of molecules in a volume V of gas at pressure P , then we know $PV = NkT$, or $P = nkT$, where $n = N/V$ is the number of molecules per unit volume. In other words, if we know the number of molecules per unit volume, we know the pressure, and vice versa: they are proportional to each other, since the temperature is constant in this problem. But the pressure is not constant, it must increase as the altitude is reduced, because it has to hold, so to speak, the weight of all the gas above it. That is the clue by which we may determine how the pressure changes with height. If we take a unit area at height h , then the vertical force from below,

40-1 The exponential atmosphere

40-2 The Boltzmann law

40-3 Evaporation of a liquid

40-4 The distribution of molecular speeds

40-5 The specific heats of gases

40-6 The failure of classical physics

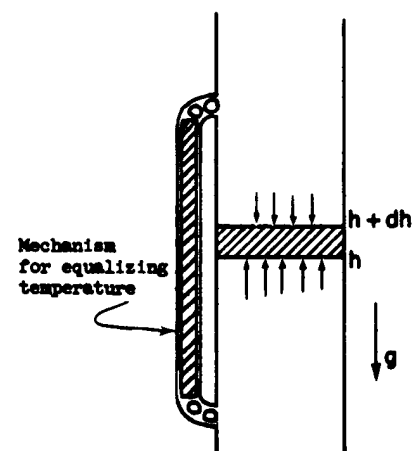


Fig. 40-1. The pressure at height h must exceed that at $h + dh$ by the weight of the intervening gas.

on this unit area, is the pressure P . The vertical force per unit area pushing down at a height $h + dh$ would be the same, in the absence of gravity, but here it is not, because the force from below must exceed the force from above by the weight of gas in the section between h and $h + dh$. Now mg is the force of gravity on each molecule, where g is the acceleration due to gravity, and $n dh$ is the total number of molecules in the unit section. So this gives us the differential equation $P_{h+dh} - P_h = dP = -mgn dh$. Since $P = nkT$, and T is constant, we can eliminate either P or n , say P , and get

$$\frac{dn}{dh} = -\frac{mg}{kT} n$$

for the differential equation, which tells us how the density goes down as we go up in energy.

We thus have an equation for the particle density n , which varies with height, but which has a derivative which is proportional to itself. Now a function which has a derivative proportional to itself is an exponential, and the solution of this differential equation is

$$n = n_0 e^{-mgh/kT}. \quad (40.1)$$

Here the constant of integration, n_0 , is obviously the density at $h = 0$ (which can be chosen anywhere), and the density goes down exponentially with height.

Note that if we have different kinds of molecules with different masses, they go down with different exponentials. The ones which were heavier would decrease with altitude faster than the light ones. Therefore we would expect that because oxygen is heavier than nitrogen, as we go higher and higher in an atmosphere with nitrogen and oxygen the proportion of nitrogen would increase. This does not really happen in our own atmosphere, at least at reasonable heights, because there is so much agitation which mixes the gases back together again. It is not an isothermal atmosphere. Nevertheless, there is a tendency for lighter materials, like hydrogen, to dominate at very great heights in the atmosphere, because the lowest masses continue to exist, while the other exponentials have all died out (Fig. 40-2).

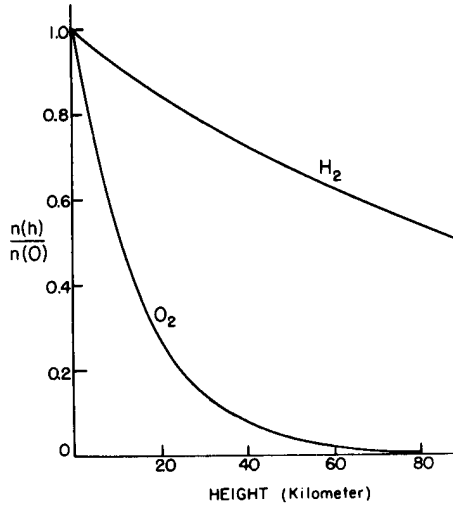


Fig. 40-2. The normalized density as a function of height in the earth's gravitational field for oxygen and for hydrogen, at constant temperature.

40-2 The Boltzmann law

Here we note the interesting fact that the numerator in the exponent of Eq. (40.1) is the *potential energy* of an atom. Therefore we can also state this particular law as: the density at any point is proportional to

$$e^{-(\text{the potential energy of each atom}/kT)}.$$

That may be an accident, i.e., may be true only for this particular case of a uniform gravitational field. However, we can show that it is a more general proposition. Suppose that there were some kind of force other than gravity acting on the molecules in a gas. For example, the molecules may be charged electrically, and may be acted on by an electric field or another charge that attracts them. Or, because of the mutual attractions of the atoms for each other, or for the wall, or for a solid, or something, there is some force of attraction which varies with position and which acts on all the molecules. Now suppose, for simplicity, that the molecules are all the same, and that the force acts on each individual one, so that the total force on a piece of gas would be simply the number of molecules times the force on each one. To avoid unnecessary complication, let us choose a coordinate system with the x -axis in the direction of the force, F .

In the same manner as before, if we take two parallel planes in the gas, separated by a distance dx , then the force on each atom, times the n atoms per cm^3 (the generalization of the previous nmg), times dx , must be balanced by the pressure change: $Fn dx = dP = kT dn$. Or, to put this law in a form which will be useful to us later,

$$F = kT \frac{d}{dx} (\ln n). \quad (40.2)$$

For the present, observe that $-F dx$ is the work we would do in taking a molecule from x to $x + dx$, and if F comes from a potential, i.e., if the work done can be represented by a potential energy at all, then this would also be the difference in the potential energy (P.E.). The negative differential of potential energy is the work done, $F dx$, and we find that $d(\ln n) = -d(\text{P.E.})/kT$, or, after integrating,

$$n = (\text{constant})e^{-\text{P.E.}/kT}. \quad (40.3)$$

Therefore what we noticed in a special case turns out to be true in general. (What if F does not come from a potential? Then (40.2) has no solution at all. Energy can be generated, or lost by the atoms running around in cyclic paths for which the work done is not zero, and no equilibrium can be maintained at all. Thermal equilibrium cannot exist if the external forces on the atoms are not conservative.) Equation (40.3), known as *Boltzmann's law*, is another of the principles of statistical mechanics: that the probability of finding molecules in a given spatial arrangement varies exponentially with the negative of the potential energy of that arrangement, divided by kT .

This, then, could tell us the distribution of molecules: Suppose that we had a positive ion in a liquid, attracting negative ions around it, how many of them would be at different distances? If the potential energy is known as a function of distance, then the proportion of them at different distances is given by this law, and so on, through many applications.

40-3 Evaporation of a liquid

In more advanced statistical mechanics one tries to solve the following important problem. Consider an assembly of molecules which attract each other, and suppose that the force between any two, say i and j , depends only on their separation r_{ij} , and can be represented as the derivative of a potential function $V(r_{ij})$. Figure 40-3 shows a form such a function might have. For $r > r_0$, the energy decreases as the molecules come together, because they attract, and then the energy increases very sharply as they come still closer together, because they repel strongly, which is characteristic of the way molecules behave, roughly speaking.

Now suppose we have a whole box full of such molecules, and we would like to know how they arrange themselves on the average. The answer is $e^{-\text{P.E.}/kT}$. The total potential energy in this case would be the sum over all the pairs, supposing that the forces are all in pairs (there may be three-body forces in more complicated things, but in electricity, for example, the potential energy is all in pairs). Then the probability for finding molecules in any particular combination of r_{ij} 's will be proportional to

$$\exp \left[- \sum_{i,j} V(r_{ij})/kT \right].$$

Now, if the temperature is very high, so that $kT \gg |V(r_0)|$, the exponent is relatively small almost everywhere, and the probability of finding a molecule is almost independent of position. Let us take the case of just two molecules: the $e^{-\text{P.E.}/kT}$ would be the probability of finding them at various mutual distances r . Clearly, where the potential goes most negative, the probability is largest, and where the potential goes toward infinity, the probability is almost zero, which occurs for very small distances. That means that for such atoms in a gas, there is no chance that they are on top of each other, since they repel so strongly. But there is a greater chance of finding them *per unit volume* at the point r_0 than at any other point. How much greater, depends on the temperature. If the temperature is very large compared with the difference in energy between $r = r_0$ and $r = \infty$, the exponential is always nearly unity. In this case, where the mean kinetic energy (about kT) greatly exceeds the potential energy, the forces do not make much difference. But as the temperature falls, the probability of finding the molecules at the preferred distance r_0 gradually increases relative to the probability of finding them elsewhere and, in fact, if kT is much less than $|V(r_0)|$, we have a relatively

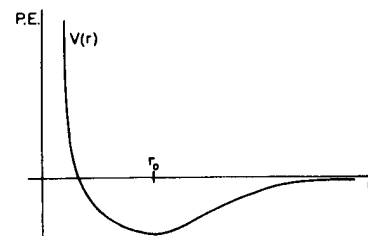


Fig. 40-3. A potential-energy function for two molecules, which depends only on their separation.

large positive exponent in that neighborhood. In other words, in a given volume they are *much more* likely to be at the distance of minimum energy than far apart. As the temperature falls, the atoms fall together, clump in lumps, and reduce to liquids, and solids, and molecules, and as you heat them up they evaporate.

The requirements for the determination of exactly how things evaporate, exactly how things should happen in a given circumstance, involve the following. First, to discover the correct molecular-force law $V(r)$, which must come from something else, quantum mechanics, say, or experiment. But, given the law of force between the molecules, to discover what a billion molecules are going to do merely consists of studying the function $e^{-\sum V_{ij}/kT}$. Surprisingly enough, since it is such a simple function and such an easy idea, given the potential, the labor is *enormously complicated*; the difficulty is the tremendous number of variables.

In spite of such difficulties, the subject is quite exciting and interesting. It is often called an example of a "many-body problem," and it really has been a very interesting thing. In that single formula must be contained all the details, for example, about the solidification of gas, or the forms of the crystals that the solid can take, and people have been trying to squeeze it out, but the mathematical difficulties are very great, not in writing the law, but in dealing with so enormous a number of variables.

That then, is the distribution of particles in space. That is the end of classical statistical mechanics, practically speaking, because if we know the forces, we can, in principle, find the distribution in space, and the distribution of velocities is something that we can work out once and for all, and is not something that is different for the different cases. The great problems are in getting particular information out of our formal solution, and that is the main subject of classical statistical mechanics.

40-4 The distribution of molecular speeds

Now we go on to discuss the distribution of velocities, because sometimes it is interesting or useful to know how many of them are moving at different speeds. In order to do that, we may make use of the facts which we discovered with regard to the gas in the atmosphere. We take it to be a perfect gas, as we have already assumed in writing the potential energy, disregarding the energy of mutual attraction of the atoms. The only potential energy that we included in our first example was gravity. We would, of course, have something more complicated if there were forces between the atoms. Thus we assume that there are no forces between the atoms and, for a moment, disregard collisions also, returning later to the justification of this. Now we saw that there are fewer molecules at the height h than there are at the height 0; according to formula (40.1), they decrease exponentially with height. How can there be fewer at greater heights? After all, do not all the molecules which are moving up at height 0 arrive at h ? No!, because some of those which are moving up at 0 are going too slowly, and cannot climb the potential hill to h . With that clue, we can calculate how many must be moving at various speeds, because from (40.1) we know *how many* are moving with less than enough speed to climb a given distance h . Those are just the ones that account for the fact that the density at h is lower than at 0.

Now let us put that idea a little more precisely: let us count how many molecules are passing from below to above the plane $h = 0$ (by calling it height = 0, we do not mean that there is a floor there; it is just a convenient label, and there is gas at negative h). These gas molecules are moving around in every direction, but some of them are moving through the plane, and at any moment a certain number per second of them are passing through the plane from below to above with different velocities. Now we note the following: if we call u the velocity which is just needed to get up to the height h (kinetic energy $mu^2/2 = mgh$), then the number of molecules per second which are passing upward through the lower plane in a vertical direction with velocity component greater than u is exactly the same as the number which pass through the upper plane with *any* upward velocity. Those molecules whose vertical velocity does not exceed u cannot get through the upper

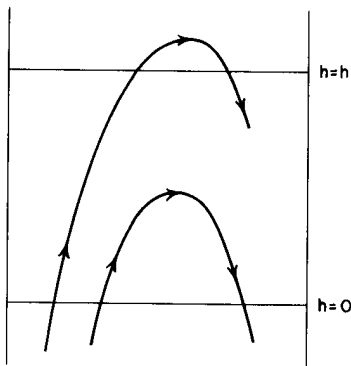


Fig. 40-4. Only those molecules moving up at $h = 0$ with sufficient velocity can arrive at height h .

plane. So therefore we see that

$$\text{Number passing } h = 0 \text{ with } v_z > u = \text{number passing } h = h \text{ with } v_z > 0.$$

But the number which pass through h with any velocity greater than 0 is less than the number which pass through the lower height with any velocity greater than 0, because the number of atoms is greater; that is all we need. We know already that the distribution of velocities is the same, after the argument we made earlier about the temperature being constant all the way through the atmosphere. So, since the velocity distributions are the same, and it is just that there are *more atoms* lower down, clearly the number $n_{>0}(h)$, passing with positive velocity at height h , and the number $n_{>0}(0)$, passing with positive velocity at height 0, are in the same ratio as the densities at the two heights, which is $e^{-mgh/kT}$. But $n_{>0}(h) = n_{>u}(0)$, and therefore we find that

$$\frac{n_{>u}(0)}{n_{>0}(0)} = e^{-mgh/kT} = e^{-mu^2/2kT},$$

since $\frac{1}{2}mu^2 = mgh$. Thus, in words, the number of molecules per unit area per second passing the height 0 with a z -component of velocity greater than u is $e^{-mu^2/2kT}$ times the total number that are passing through the plane with velocity greater than zero.

Now this is not only true at the arbitrarily chosen height 0, but of course it is true at any other height, and thus the distributions of velocities *are* all the same! (The final statement does not involve the height h , which appeared only in the intermediate argument.) The result is a general proposition that gives us the distribution of velocities. It tells us that if we drill a little hole in the side of a gas pipe, a very tiny hole, so that the collisions are few and far between, i.e., are farther apart than the diameter of the hole, then the particles which are coming out will have different velocities, but the fraction of particles which come out at a velocity greater than u is $e^{-mu^2/2kT}$.

Now we return to the question about the neglect of collisions: Why does it not make any difference? We could have pursued the same argument, not with a finite height h , but with an infinitesimal height h , which is so small that there would be no room for collisions between 0 and h . But that was not necessary: the argument is evidently based on an analysis of the energies involved, the conservation of energy, and in the collisions that occur there is an exchange of energies among the molecules. However, we do not really care whether we follow the same molecule if energy is merely exchanged with another molecule. So it turns out that even if the problem is analyzed more carefully (and it is more difficult, naturally, to do a rigorous job), it still makes no difference in the result.

It is interesting that the velocity distribution we have found is just

$$n_{>u} \propto e^{-\text{kinetic energy}/kT}. \quad (40.4)$$

This way of describing the distribution of velocities, by giving the number of molecules that pass a given area with a certain minimum z -component, is not the most convenient way of giving the velocity distribution. For instance, inside the gas, one more often wants to know how many molecules are moving with a z -component of velocity between two given values, and that, of course, is not directly given by Eq. (40.4). We would like to state our result in the more conventional form, even though what we already have written is quite general. *Note that it is not possible to say that any molecule has exactly some stated velocity; none of them has a velocity exactly equal to 1.7962899173 meters per second.* So in order to make a meaningful statement, we have to ask how many are to be found in some *range* of velocities. We have to say how many have velocities between 1.796 and 1.797, and so on. On mathematical terms, let $f(u) du$ be the fraction of all the molecules which have velocities between u and $u + du$ or, what is the same thing (if du is infinitesimal), all that have a velocity u with a range du . Figure 40-5 shows a possible form for the function $f(u)$, and the shaded part, of width du and mean height $f(u)$, represents this fraction $f(u) du$. That is, the ratio of the shaded

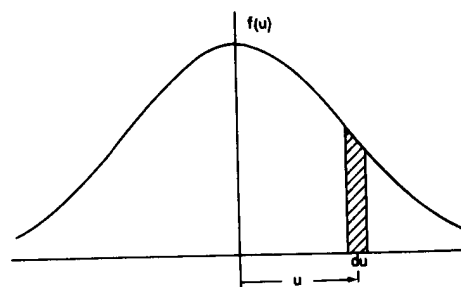


Fig. 40-5. A velocity distribution function. The shaded area is $f(u) du$, the fraction of particles having velocities within a range du about u .

area to the total area of the curve is the relative proportion of molecules with velocity u within du . If we define $f(u)$ so that the fraction having a velocity in this range is given directly by the shaded area, then the total area must be 100 percent of them, that is,

$$\int_{-\infty}^{\infty} f(u) du = 1. \quad (40.5)$$

Now we have only to get this distribution by comparing it with the theorem we derived before. First we ask, what is the number of molecules passing through an area per second with a velocity greater than u , expressed in terms of $f(u)$? At first we might think it is merely the integral of $\int_u^{\infty} f(u) du$, but it is not, because we want the number that are passing the area *per second*. The faster ones pass more often, so to speak, than the slower ones, and in order to express how many pass, you have to multiply by the velocity. (We discussed that in the previous chapter when we talked about the number of collisions.) In a given time t the total number which pass through the surface is all of those which have been able to arrive at the surface, and the number which arrive come from a distance ut . So the number of molecules which arrive is not simply the number which are there, but the number that are there per unit volume, multiplied by the distance that they sweep through in racing for the area through which they are supposed to go, and that distance is proportional to u . Thus we need the integral of u times $f(u) du$, an infinite integral with a lower limit u , and this must be the same as we found before, namely $e^{-mu^2/2kT}$, with a proportionality constant which we will get later:

$$\int_u^{\infty} uf(u) du = \text{const} \cdot e^{-mu^2/2kT}. \quad (40.6)$$

Now if we differentiate the integral with respect to u , we get the thing that is inside the integral, i.e., the integrand (with a minus sign, since u is the lower limit), and if we differentiate the other side, we get u times the same exponential (and some constants). The u 's cancel and we find

$$f(u) du = Ce^{-mu^2/2kT} du. \quad (40.7)$$

We retain the du on both sides as a reminder that it is a *distribution*, and it tells what the proportion is for velocity between u and $u + du$.

The constant C must be so determined that the integral is unity, according to Eq. (40.5). Now we can prove* that

$$\int_{-\infty}^{\infty} e^{-x^2} dx = \sqrt{\pi}.$$

Using this fact, it is easy to find that $C = \sqrt{m/2\pi kT}$.

Since velocity and momentum are proportional, we may say that the distribution of momenta is also proportional to $e^{-K.E./kT}$ per unit momentum range. It turns out that this theorem is true in relativity too, if it is in terms of momentum, while if it is in velocity it is not, so it is best to learn it in momentum instead of in velocity:

$$f(p) dp = Ce^{-K.E./kT} dp. \quad (40.8)$$

So we find that the probabilities of different conditions of energy, kinetic and potential, are both given by $e^{-\text{energy}/kT}$, a very easy thing to remember and a rather beautiful proposition.

* To get the value of the integral, let

$$I = \int_{-\infty}^{\infty} e^{-x^2} dx.$$

Then

$$I^2 = \int_{-\infty}^{\infty} e^{-x^2} dx \cdot \int_{-\infty}^{\infty} e^{-y^2} dy = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} e^{-(x^2+y^2)} dy dx,$$

which is a double integral over the whole xy -plane. But this can also be written in polar coordinates as

$$I^2 = \int_0^{\infty} e^{-r^2} \cdot 2\pi r dr = \pi \int_0^{\infty} e^{-t} dt = \pi.$$

So far we have, of course, only the distribution of the velocities “vertically.” We might want to ask, what is the probability that a molecule is moving in another direction? Of course these distributions are connected, and one can obtain the complete distribution from the one we have, because the complete distribution depends only on the square of the magnitude of the velocity, not upon the z -component. It must be something that is independent of direction, and there is only one function involved, the probability of different magnitudes. We have the distribution of the z -component, and therefore we can get the distribution of the other components from it. The result is that the probability is still proportional to $e^{-K.E./kT}$, but now the kinetic energy involves three parts, $mv_x^2/2$, $mv_y^2/2$, and $mv_z^2/2$, summed in the exponent. Or we can write it as a product:

$$f(v_x, v_y, v_z) dv_x dv_y dv_z \propto e^{-mv_x^2/2kT} \cdot e^{-mv_y^2/2kT} \cdot e^{-mv_z^2/2kT} dv_x dv_y dv_z. \quad (40.9)$$

You can see that this formula must be right because, first, it is a function only of v^2 , as required, and second, the probabilities of various values of v_z obtained by integrating over all v_x and v_y is just (40.7). But this one function (40.9) can do both those things!

40-5 The specific heats of gases

Now we shall look at some ways to test the theory, and to see how successful is the classical theory of gases. We saw earlier that if U is the internal energy of N molecules, then $PV = NkT = (\gamma - 1)U$ holds, sometimes, for some gases, maybe. If it is a monatomic gas, we know this is also equal to $\frac{2}{3}$ of the kinetic energy of the center-of-mass motion of the atoms. If it is a monatomic gas, then the kinetic energy is equal to the internal energy, and therefore $\gamma - 1 = \frac{2}{3}$. But suppose it is, say, a more complicated molecule, that can spin and vibrate, and let us suppose (it turns out to be true according to classical mechanics) that the energies of the internal motions are also proportional to kT . Then at a given temperature, in addition to kinetic energy kT , it has internal vibrational or rotational energy. So the total U includes not just the internal kinetic energy, but also the rotational energy, and we get a different value of γ . Technically, the best way to measure γ is by measuring the specific heat, which is the change in energy with temperature. We will return to that approach later. For our present purposes, we may suppose γ is found experimentally from the PV^γ curve for adiabatic compression.

Let us make a calculation of γ for some cases. First, for a monatomic gas U is the total energy, the same as the kinetic energy, and we know already that γ should be $\frac{5}{3}$. For a diatomic gas, we may take, as an example, oxygen, hydrogen iodide, hydrogen, etc., and suppose that the diatomic gas can be represented as two atoms held together by some kind of force like the one of Fig. 40-3. We may also suppose, and it turns out to be quite true, that at the temperatures that are of interest for the diatomic gas, the pairs of atoms tend strongly to be separated by r_0 , the distance of potential minimum. If this were not true, if the probability were not strongly varying enough to make the great majority sit near the bottom, we would have to remember that oxygen gas is a mixture of O_2 and single oxygen atoms in a nontrivial ratio. We know that there are, in fact, very few single oxygen atoms, which means that the potential energy minimum is very much greater in magnitude than kT , as we have seen. Since they are clustered strongly around r_0 , the only part of the curve that is needed is the part near the minimum, which may be approximated by a parabola. A parabolic potential implies a harmonic oscillator, and in fact, to an excellent approximation, the oxygen molecule can be represented as two atoms connected by a spring.

Now what is the total energy of this molecule at temperature T ? We know that for each of the two atoms, each of the kinetic energies should be $\frac{3}{2}kT$, so the kinetic energy of both of them is $\frac{3}{2}kT + \frac{3}{2}kT$. We can also put this in a different way: the same $\frac{3}{2}$ plus $\frac{3}{2}$ can also be looked at as kinetic energy of the center of mass

Table 40-1

Values of the specific heat ratio, γ ,
for various gases

Gas	T ($^{\circ}\text{C}$)	γ
He	-180	1.660
Kr	19	1.68
Ar	15	1.668
H ₂	100	1.404
O ₂	100	1.399
HI	100	1.40
Br ₂	300	1.32
I ₂	185	1.30
NH ₃	15	1.310
C ₂ H ₆	15	1.22

($\frac{3}{2}$), kinetic energy of rotation ($\frac{2}{2}$), and kinetic energy of vibration ($\frac{1}{2}$). We know that the kinetic energy of vibration is $\frac{1}{2}$, since there is just one dimension involved and each degree of freedom has $\frac{1}{2}kT$. Regarding the rotation, it can turn about either of two axes, so there are two independent motions. We assume that the atoms are some kind of points, and cannot spin about the line joining them; this is something to bear in mind, because if we get a disagreement, maybe that is where the trouble is. But we have one more thing, which is the *potential* energy of vibration; how much is that? In a harmonic oscillator the average kinetic energy and average potential energy are equal, and therefore the potential energy of vibration is $\frac{1}{2}kT$, also. The grand total of energy is $U = \frac{7}{2}kT$, or kT is $\frac{2}{7}U$ per atom. That means, then, that γ is $\frac{9}{7}$ instead of $\frac{5}{3}$, i.e., $\gamma = 1.286$.

We may compare these numbers with the relevant measured values shown in Table 40-1. Looking first at helium, which is a monatomic gas, we find very nearly $\frac{5}{3}$, and the error is probably experimental, although at such a low temperature there may be some forces between the atoms. Krypton and argon, both monatomic, agree also within the accuracy of the experiment.

We turn to the diatomic gases and find hydrogen with 1.404, which does not agree with the theory, 1.286. Oxygen, 1.399, is very similar, but again not in agreement. Hydrogen iodide again is similar at 1.40. It begins to look as though the right answer is 1.40, but it is not, because if we look further at bromine we see 1.32, and at iodine we see 1.30. Since 1.30 is reasonably close to 1.286, iodine may be said to agree rather well, but oxygen is far off. So here we have a dilemma. We have it right for one molecule, we do not have it right for another molecule, and we may need to be pretty ingenious in order to explain both.

Let us look further at a still more complicated molecule with large numbers of parts, for example, C₂H₆, which is ethane. It has eight different atoms, and they are all vibrating and rotating in various combinations, so the total amount of internal energy must be an enormous number of kT 's, at least $12kT$ for kinetic energy alone, and $\gamma - 1$ must be very close to zero, or γ almost exactly 1. In fact, it is lower, but 1.22 is not so much lower, and is higher than the $1\frac{1}{2}$ calculated from the kinetic energy alone, and it is just not understandable!

Furthermore, the whole mystery is deep, because the diatomic molecule cannot be made rigid by a limit. Even if we made the couplings stiffer indefinitely, although it might not vibrate much, it would nevertheless keep vibrating. The vibrational energy inside is still kT , since it does not depend on the strength of the coupling. But if we could imagine *absolute* rigidity, stopping all vibration to eliminate a variable, then we would get $U = \frac{5}{2}kT$ and $\gamma = 1.40$ for the diatomic case. This looks good for H₂ or O₂. On the other hand, we would still have problems, because γ for either hydrogen or oxygen varies with temperature! From the measured values shown in Fig. 40-6, we see that for H₂, γ varies from about 1.6 at -185°C to 1.3 at 2000°C . The variation is more substantial in the case of hydrogen than for oxygen, but nevertheless, even in oxygen, γ tends definitely to go up as we go down in temperature.

40-6 The failure of classical physics

So, all in all, we might say that we have some difficulty. We might try some force law other than a spring, but it turns out that anything else will only make γ higher. If we include more forms of energy, γ approaches unity more closely, contradicting the facts. All the classical theoretical things that one can think of will only make it worse. The fact is that there are electrons in each atom, and we know from their spectra that there are internal motions; each of the electrons should have at least $\frac{1}{2}kT$ of kinetic energy, and something for the potential energy, so when these are added in, γ gets still smaller. It is ridiculous. It is wrong.

The first great paper on the dynamical theory of gases was by Maxwell in 1859. On the basis of ideas we have been discussing, he was able accurately to explain a great many known relations, such as Boyle's law, the diffusion theory, the viscosity of gases, and things we shall talk about in the next chapter. He listed all these great successes in a final summary, and at the end he said, "Finally, by

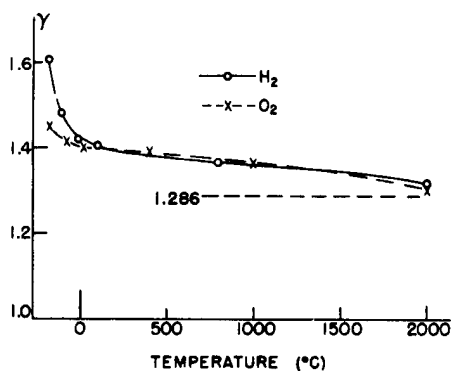


Fig. 40-6. Experimental values of γ as a function of temperature for hydrogen and oxygen. Classical theory predicts $\gamma = 1.286$, independent of temperature.

establishing a necessary relation between the motions of translation and rotation (he is talking about the $\frac{1}{2}kT$ theorem) of all particles not spherical, we proved that a system of such particles could not possibly satisfy the known relation between the two specific heats." He is referring to γ (which we shall see later is related to two ways of measuring specific heat), and he says we know we cannot get the right answer.

Ten years later, in a lecture, he said, "I have now put before you what I consider to be the greatest difficulty yet encountered by the molecular theory." These words represent the first discovery that the laws of classical physics were wrong. This was the first indication that there was something fundamentally impossible, because a rigorously proved theorem did not agree with experiment. About 1890, Jeans was to talk about this puzzle again. One often hears it said that physicists at the latter part of the nineteenth century thought they knew all the significant physical laws and that all they had to do was to calculate more decimal places. Someone may have said that once, and others copied it. But a thorough reading of the literature of the time shows they were all worrying about something. Jeans said about this puzzle that it is a very mysterious phenomenon, and it seems as though as the temperature falls, certain kinds of motions "freeze out."

If we could assume that the vibrational motion, say, did not exist at low temperature and did exist at high temperature, then we could imagine that a gas might exist at a temperature sufficiently low that vibrational motion does not occur, so $\gamma = 1.40$, or a higher temperature at which it begins to come in, so γ falls. The same might be argued for the rotation. If we can eliminate the rotation, say it "freezes out" at sufficiently low temperature, then we can understand the fact that the γ of hydrogen approaches 1.66 as we go down in temperature. How can we understand such a phenomenon? Of course that these motions "freeze out" cannot be understood by classical mechanics. It was only understood when quantum mechanics was discovered.

Without proof, we may state the results for statistical mechanics of the quantum-mechanical theory. We recall that according to quantum mechanics, a system which is bound by a potential, for the vibrations, for example, will have a discrete set of energy levels, i.e., states of different energy. Now the question is: how is statistical mechanics to be modified according to quantum-mechanical theory? It turns out, interestingly enough, that although most problems are more difficult in quantum mechanics than in classical mechanics, problems in statistical mechanics are much easier in quantum theory! The simple result we have in classical mechanics, that $n = n_0 e^{-\text{energy}/kT}$, becomes the following very important theorem: If the energies of the set of molecular states are called, say, $E_0, E_1, E_2, \dots, E_i, \dots$, then in thermal equilibrium the probability of finding a molecule in the particular state of having energy E_i is proportional to $e^{-E_i/kT}$. That gives the probability of being in various states. In other words, the relative chance, the probability, of being in state E_1 relative to the chance of being in state E_0 , is

$$\frac{P_1}{P_0} = \frac{e^{-E_1/kT}}{e^{-E_0/kT}}, \quad (40.10)$$

which, of course, is the same as

$$n_1 = n_0 e^{-(E_1 - E_0)/kT}, \quad (40.11)$$

since $P_1 = n_1/N$ and $P_0 = n_0/N$. So it is less likely to be in a higher energy state than in a lower one. The ratio of the number of atoms in the upper state to the number in the lower state is e raised to the power (minus the energy difference, over kT)—a very simple proposition.

Now it turns out that for a harmonic oscillator the energy levels are evenly spaced. Calling the lowest energy $E_0 = 0$ (it actually is not zero, it is a little different, but it does not matter if we shift all energies by a constant), the first one is then $E_1 = \hbar\omega$, and the second one is $2\hbar\omega$, and the third one is $3\hbar\omega$, and so on.

Now let us see what happens. We suppose we are studying the vibrations of a diatomic molecule, which we approximate as a harmonic oscillator. Let us ask

what is the relative chance of finding a molecule in state E_1 instead of in state E_0 . The answer is that the chance of finding it in state E_1 , relative to that of finding it in state E_0 , goes down as $e^{-\hbar\omega/kT}$. Now suppose that kT is much less than $\hbar\omega$, and we have a low-temperature circumstance. Then the probability of its being in state E_1 is extremely small. Practically all the atoms are in state E_0 . If we change the temperature but still keep it very small, then the chance of its being in state $E_1 = \hbar\omega$ remains infinitesimal—the energy of the oscillator remains nearly zero; it does not change with temperature so long as the temperature is much less than $\hbar\omega$. All oscillators are in the bottom state, and their motion is effectively “frozen”; *there is no contribution of it to the specific heat*. We can judge, then, from Table 40-1, that at 100°C, which is 373 degrees absolute, kT is much less than the vibrational energy in the oxygen or hydrogen molecules, but not so in the iodine molecule. The reason for the difference is that an iodine atom is very heavy, compared with hydrogen, and although the forces may be comparable in iodine and hydrogen, the iodine molecule is so heavy that the natural frequency of vibration is very low compared with the natural frequency of hydrogen. With $\hbar\omega$ higher than kT at room temperature for hydrogen, but lower for iodine, only the latter, iodine, exhibits the classical vibrational energy. As we increase the temperature of a gas, starting from a very low value of T , with the molecules almost all in their lowest state, they gradually begin to have an appreciable probability to be in the second state, and then in the next state, and so on. When the probability is appreciable for many states, the behavior of the gas approaches that given by classical physics, because the quantized states become nearly indistinguishable from a continuum of energies, and the system can have almost any energy. Thus, as the temperature rises, we should again get the results of classical physics, as indeed seems to be the case in Fig. 40-6. It is possible to show in the same way that the rotational states of atoms are also quantized, but the states are so much closer together that in ordinary circumstances kT is bigger than the spacing. Then many levels are excited, and the rotational kinetic energy in the system participates in the classical way. The one example where this is not quite true at room temperature is for hydrogen.

This is the first time that we have really deduced, by comparison with experiment, that there was something wrong with classical physics, and we have looked for a resolution of the difficulty in quantum mechanics in much the same way as it was done originally. It took 30 or 40 years before the next difficulty was discovered, and that had to do again with statistical mechanics, but this time the mechanics of a photon gas. That problem was solved by Planck, in the early years of this century.

The Brownian Movement

41-1 Equipartition of energy

The Brownian movement was discovered in 1827 by Robert Brown, a botanist. While he was studying microscopic life, he noticed little particles of plant pollens jiggling around in the liquid he was looking at in the microscope, and he was wise enough to realize that these were not living, but were just little pieces of dirt moving around in the water. In fact he helped to demonstrate that this had nothing to do with life by getting from the ground an old piece of quartz in which there was some water trapped. It must have been trapped for millions and millions of years, but inside he could see the same motion. What one sees is that very tiny particles are jiggling all the time.

This was later proved to be one of the effects of *molecular motion*, and we can understand it qualitatively by thinking of a great push ball on a playing field, seen from a great distance, with a lot of people underneath, all pushing the ball in various directions. We cannot see the people because we imagine that we are too far away, but we can see the ball, and we notice that it moves around rather irregularly. We also know, from the theorems that we have discussed in previous chapters, that the mean kinetic energy of a small particle suspended in a liquid or a gas will be $\frac{3}{2}kT$ even though it is very heavy compared with a molecule. If it is very heavy, that means that the speeds are relatively slow, but it turns out, actually, that the speed is not really so slow. In fact, we cannot see the speed of such a particle very easily because although the mean kinetic energy is $\frac{3}{2}kT$, which represents a speed of a millimeter or so per second for an object a micron or two in diameter, this is very hard to see even in a microscope, because the particle continuously reverses its direction and does not get anywhere. How far it does get we will discuss at the end of the present chapter. This problem was first solved by Einstein at the beginning of the present century.

Incidentally, when we say that the mean kinetic energy of this particle is $\frac{3}{2}kT$, we claim to have derived this result from the kinetic theory, that is, from Newton's laws. We shall find that we can derive all kinds of things—marvelous things—from the kinetic theory, and it is most interesting that we can apparently get so much from so little. Of course we do not mean that Newton's laws are “little”—they are enough to do it, really—what we mean is that *we* did not do very much. How do we get so much out? The answer is that we have been perpetually making a certain important assumption, which is that if a given system is in thermal equilibrium at some temperature, it will also be in thermal equilibrium with *anything else* at the same temperature. For instance, if we wanted to see how a particle would move if it was really colliding with water, we could imagine that there was a gas present, composed of another kind of particle, little fine pellets that (we suppose) do not interact with water, but only hit the particle with “hard” collisions. Suppose the particle has a prong sticking out of it; all our pellets have to do is hit the prong. We know all about this imaginary gas of pellets at temperature T —it is an ideal gas. Water is complicated, but an ideal gas is simple. Now, *our particle has to be in equilibrium with the gas of pellets*. Therefore, the mean motion of the particle must be what we get for gaseous collisions, because if it were not moving at the right speed relative to the water but, say, was moving faster, that would mean that the pellets would pick up energy from it and get hotter than the water. But we had started them at the same temperature, and we assume that if a thing is once in equilibrium, it stays in equilibrium—parts of it do not get hotter and other parts colder, spontaneously.

41-1 Equipartition of energy

41-2 Thermal equilibrium of radiation

41-3 Equipartition and the quantum oscillator

41-4 The random walk

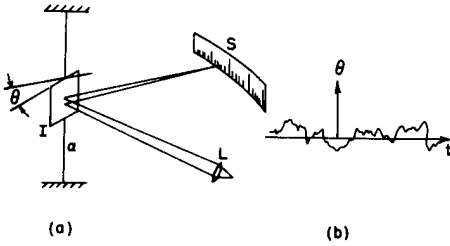


Fig. 41-1. (a) A sensitive light-beam galvanometer. Light from a source L is reflected from a small mirror onto a scale. (b) A schematic record of the reading of the scale as a function of the time.

This proposition is true and can be proved from the laws of mechanics, but the proof is very complicated and can be established only by using advanced mechanics. It is much easier to prove in quantum mechanics than it is in classical mechanics. It was proved first by Boltzmann, but for now we simply take it to be true, and then we can argue that our particle has to have $\frac{3}{2}kT$ of energy if it is hit with artificial pellets, so it also must have $\frac{3}{2}kT$ when it is being hit with water at the same temperature and we take away the pellets; so it is $\frac{3}{2}kT$. It is a strange line of argument, but perfectly valid.

In addition to the motion of colloidal particles for which the Brownian movement was first discovered, there are a number of other phenomena, both in the laboratory and in other situations, where one can see Brownian movement. If we are trying to build the most delicate possible equipment, say a very small mirror on a thin quartz fiber for a very sensitive ballistic galvanometer (Fig. 41-1), the mirror does not stay put, but jiggles all the time—all the time—so that when we shine a light on it and look at the position of the spot, we do not have a perfect instrument because the mirror is always jiggling. Why? Because the average kinetic energy of rotation of this mirror has to be, on the average, $\frac{1}{2}kT$.

What is the mean-square angle over which the mirror will wobble? Suppose we find the natural vibration period of the mirror by tapping on one side and seeing how long it takes to oscillate back and forth, and we also know the moment of inertia, I . We know the formula for the kinetic energy of rotation—it is given by Eq. (19.8): $T = \frac{1}{2}I\omega^2$. That is the kinetic energy, and the potential energy that goes with it will be proportional to the square of the angle—it is $V = \frac{1}{2}\alpha\theta^2$. But, if we know the period t_0 and calculate from that the natural frequency $\omega_0 = 2\pi/t_0$, then the potential energy is $V = \frac{1}{2}I\omega_0^2\theta^2$. Now we know that the average kinetic energy is $\frac{1}{2}kT$, but since it is a harmonic oscillator the average potential energy is also $\frac{1}{2}kT$. Thus

$$\frac{1}{2}I\omega_0^2\langle\theta^2\rangle = \frac{1}{2}kT,$$

or

$$\langle\theta^2\rangle = kT/I\omega_0^2. \quad (41.1)$$

In this way we can calculate the oscillations of a galvanometer mirror, and thereby find what the limitations of our instrument will be. If we want to have smaller oscillations, we have to cool the mirror. An interesting question is, *where* to cool it. This depends upon where it is getting its “kicks” from. If it is through the fiber, we cool it at the top—if the mirror is surrounded by a gas and is getting hit mostly by collisions in the gas, it is better to cool the gas. As a matter of fact, if we know where the *damping* of the oscillations comes from, it turns out that that is always the *source* of the fluctuations also, a point which we will come back to.

The same thing works, amazingly enough, in *electrical circuits*. Suppose that we are building a very sensitive, accurate amplifier for a definite frequency and have a resonant circuit (Fig. 41-2) in the input so as to make it very sensitive to this certain frequency, like a radio receiver, but a really good one. Suppose we wish to go down to the very lowest limit of things, so we take the voltage, say off the inductance, and send it into the rest of the amplifier. Of course, in any circuit like this, there is a certain amount of loss. It is not a perfect resonant circuit, but it is a very good one and there is a little resistance, say (we put the resistor in so we can see it, but it is supposed to be small). Now we would like to find out: How much does the voltage across the inductance fluctuate? *Answer*: We know that $\frac{1}{2}LI^2$ is the “kinetic energy”—the energy associated with a coil in a resonant circuit (Chapter 25). Therefore the mean value of $\frac{1}{2}LI^2$ is equal to $\frac{1}{2}kT$ —that tells us what the rms current is and we can find out what the rms voltage is from the rms current. For if we want the voltage across the inductance the formula is $\hat{V}_L = i\omega L\hat{I}$, and the mean absolute square voltage on the inductance is $\langle V_L^2 \rangle = L^2\omega_0^2\langle I^2 \rangle$, and putting in $\frac{1}{2}L\langle I^2 \rangle = \frac{1}{2}kT$, we obtain

$$\langle V_L^2 \rangle = L\omega_0^2 kT. \quad (41.2)$$

So now we can design circuits and tell when we are going to get what is called *Johnson noise*, the noise associated with thermal fluctuations!

Where do the fluctuations come from this time? They come again from the *resistor*—they come from the fact that the electrons in the resistor are jiggling around because they are in thermal equilibrium with the matter in the resistor, and they make fluctuations in the density of electrons. They thus make tiny electric fields which drive the resonant circuit.

Electrical engineers represent the answer in another way. Physically, the resistor is effectively the source of noise. However, we may replace the real circuit having an honest, true physical resistor which is making noise, by an artificial circuit which contains a little generator that is going to *represent* the noise, and now the resistor is otherwise ideal—no noise comes from it. All the noise is in the artificial generator. And so if we knew the characteristics of the noise generated by a resistor, if we had the formula for that, then we could calculate what the circuit is going to do in response to that noise. So, we need a formula for the noise fluctuations. Now the noise that is generated by the resistor is at all frequencies, since the resistor by itself is not resonant. Of course the resonant circuit only “listens” to the part that is near the right frequency, but the resistor has many different frequencies in it. We may describe how strong the generator is, as follows: The mean power that the resistor would absorb if it were connected directly across the noise generator would be $\langle E^2 \rangle / R$, if E were the voltage from the generator. But we would like to know in more detail how much power there is at every frequency. There is very little power in any one frequency; it is a distribution. Let $P(\omega) d\omega$ be the power that the generator would deliver in the frequency range $d\omega$ into the very same resistor. Then we can prove (we shall prove it for another case, but the mathematics is exactly the same) that the power comes out

$$P(\omega) d\omega = (2/\pi)kT d\omega, \quad (41.3)$$

and is *independent of the resistance* when put this way.

41-2 Thermal equilibrium of radiation

Now we go on to consider a still more advanced and interesting proposition that is as follows. Suppose we have a charged oscillator like those we were talking about when we were discussing light, let us say an electron oscillating up and down in an atom. If it oscillates up and down, it radiates light. Now suppose that this oscillator is in a very thin gas of other atoms, and that from time to time the atoms collide with it. Then in equilibrium, after a long time, this oscillator will pick up energy such that its kinetic energy of oscillation is $\frac{1}{2}kT$, and since it is a harmonic oscillator, its entire energy of motion will become kT . That is, of course, a wrong description so far, because the oscillator carries *electric charge*, and if it has an energy kT it is shaking up and down and *radiating light*. Therefore it is impossible to have equilibrium of real matter alone without the charges in it emitting light, and as light is emitted, energy flows away, the oscillator loses its kT as time goes on, and thus the whole gas which is colliding with the oscillator gradually cools off. And that is, of course, the way a hot stove cools on a cold night by radiating the light into the sky, because the atoms are jiggling their charge and they continually radiate, and slowly, because of this radiation, the jiggling motion slows down.

On the other hand, if we enclose the whole thing in a box so that the light does not go away to infinity, then we *can* eventually get thermal equilibrium. We may either put the gas in a box where we can say that there are other radiators in the box walls sending light back or, to take a nicer example, we may suppose the box has mirror walls. It is easier to think about that case. Thus we assume that all the radiation that goes out from the oscillator keeps running around in the box. Then, of course, it is true that the oscillator starts to radiate, but pretty soon it can maintain its kT of kinetic energy in spite of the fact that it is radiating, because it is being illuminated, we may say, by its own light reflected from the walls of the box. That is, after a while there is a great deal of light rushing around in the box, and although the oscillator is radiating some, the light comes back and returns some of the energy that was radiated.

We shall now determine how much light there must be in such a box at temperature T in order that the shining of the light on this oscillator will generate just enough energy to account for the light it radiated.

Let the gas atoms be very few and far between, so that we have an ideal oscillator with no resistance except radiation resistance. Then we consider that at thermal equilibrium the oscillator is doing two things at the same time. First, it has a mean energy kT , and we calculate how much radiation it emits. Second, this radiation should be exactly the amount that would result because of the fact that the light shining on the oscillator is scattered. Since there is nowhere else the energy can go, this effective radiation is really just scattered light from the light that is in there.

Thus we first calculate the energy that is radiated by the oscillator per second, if the oscillator has a certain energy. (We borrow from Chapter 32 on radiation resistance a number of equations without going back over their derivation.) The energy radiated per radian divided by the energy of the oscillator is called $1/Q$ (Eq. 32.8): $1/Q = (dW/dt)/\omega_0 W$. Using the quantity γ , the damping constant, this can also be written as $1/Q = \gamma/\omega_0$, where ω_0 is the natural frequency of the oscillator—if gamma is very small, Q is very large. The energy radiated per second is then

$$\frac{dW}{dt} = \frac{\omega_0 W}{Q} = \frac{\omega_0 W \gamma}{\omega_0} = \gamma W. \quad (41.4)$$

The energy radiated per second is thus simply gamma times the energy of the oscillator. Now the oscillator should have an average energy kT , so we see that gamma kT is the average amount of energy radiated per second:

$$\langle dW/dt \rangle = \gamma kT. \quad (41.5)$$

Now we only have to know what gamma is. Gamma is easily found from Eq. (32.12). It is

$$\gamma = \frac{\omega_0}{Q} = \frac{2}{3} \frac{r_0 \omega_0^2}{c}, \quad (41.6)$$

where $r_0 = e^2/mc^2$ is the classical electron radius, and we have set $\lambda = 2\pi c/\omega_0$.

Our final result for the average rate of radiation of light near the frequency ω_0 is therefore

$$\frac{dW}{dt} = \frac{2}{3} \frac{r_0 \omega_0^2 kT}{c}. \quad (41.7)$$

Next we ask how much light must be shining on the oscillator. It must be enough that the energy absorbed from the light (and thereupon scattered) is just exactly this much. In other words, the emitted light is accounted for as *scattered* light from the light that is shining on the oscillator in the cavity. So we must now calculate how much light is scattered from the oscillator if there is a certain amount—unknown—of radiation incident on it. Let $I(\omega) d\omega$ be the amount of light energy there is at the frequency ω , within a certain range $d\omega$ (because there is no light at *exactly* a certain frequency; it is spread all over the spectrum). So $I(\omega)$ is a certain *spectral distribution* which we are now going to find—it is the color of a furnace at temperature T that we see when we open the door and look in the hole. Now how much light is absorbed? We worked out the amount of radiation absorbed from a given incident light beam, and we calculated it in terms of a *cross section*. It is just as though we said that all of the light that falls on a certain cross section is absorbed. So the total amount that is re-radiated (scattered) is the incident intensity $I(\omega) d\omega$ multiplied by the cross section σ .

The formula for the cross section which we derived (Eq. 31.19) did not have the damping included. It is not hard to go through the derivation again and put in the resistance term which we neglected. If we do that, and calculate the cross section the same way, we get

$$\sigma_s = \frac{8\pi r_0^2}{3} \left(\frac{\omega^4}{(\omega^2 - \omega_0^2)^2 + \gamma^2 \omega^2} \right). \quad (41.8)$$

Now, as a function of frequency, σ_s is of significant size only for ω very near to the natural frequency ω_0 . (Remember that the Q for a radiating oscillator is about 10^8 .) The oscillator scatters very strongly when ω is equal to ω_0 , and very weakly for other values of ω . Therefore we can replace ω by ω_0 and $\omega^2 - \omega_0^2$ by $2\omega_0(\omega - \omega_0)$, and we get

$$\sigma_s = \frac{2\pi r_0^2 \omega_0^2}{3[(\omega - \omega_0)^2 + \gamma^2/4]} \quad (41.9)$$

Now the whole curve is localized near $\omega = \omega_0$. (We do not really have to make any approximations, but it is much easier to do the integrals if we simplify the equation a bit.) Now we multiply the intensity in a given frequency range by the cross section of scattering, to get the amount of energy scattered in the range $d\omega$. The *total* energy scattered is then the integral of this for all ω . Thus

$$\begin{aligned} \frac{dW_s}{dt} &= \int_0^\infty I(\omega) \sigma_s(\omega) d\omega \\ &= \int_0^\infty \frac{2\pi r_0^2 \omega_0^2 I(\omega) d\omega}{3[(\omega - \omega_0)^2 + \gamma^2/4]} \end{aligned} \quad (41.10)$$

Now we set $dW_s/dt = 3\gamma kT$. Why *three*? Because when we made our analysis of the cross section in Chapter 32, we assumed that the polarization was such that the light could drive the oscillator. If we had used an oscillator which could move only in one direction, and the light, say, was polarized in the wrong way, it would not give any scattering. So we must either average the cross section of an oscillator which can go only in one direction, over all directions of incidence and polarization of the light or, more easily, we can imagine an oscillator which will follow the field no matter which way the field is pointing. Such an oscillator, which can oscillate equally in three directions, would have $3kT$ average energy because there are 3 degrees of freedom in that oscillator. So we should use $3\gamma kT$ because of the 3 degrees of freedom.

Now we have to do the integral. Let us suppose that the unknown spectral distribution $I(\omega)$ of the light is a smooth curve and does not vary very much across the very narrow frequency region where σ_s is peaked (Fig. 41-3). Then the only significant contribution comes when ω is very close to ω_0 , within an amount gamma, which is very small. So therefore, although $I(\omega)$ may be an unknown and complicated function, the only place where it is important is near $\omega = \omega_0$, and there we may replace the smooth curve by a flat one—a “constant”—at the same height. In other words, we simply take $I(\omega)$ outside the integral sign and call it $I(\omega_0)$. We may also take the rest of the constants out in front of the integral, and what we have left is

$$\frac{2}{3} \pi r_0^2 \omega_0^2 I(\omega_0) \int_0^\infty \frac{d\omega}{(\omega - \omega_0)^2 + \gamma^2/4} = 3\gamma kT. \quad (41.11)$$

Now, the integral should go from 0 to ∞ , but 0 is so far from ω_0 that the curve is all finished by that time, so we go instead to minus ∞ —it makes no difference and it is much easier to do the integral. The integral is an inverse tangent function of the form $\int dx/(x^2 + a^2)$. If we look it up in a book we see that it is equal to π/a . So what it comes to for our case is $2\pi/\gamma$. Therefore we get, with some rearranging,

$$I(\omega_0) = \frac{9\gamma^2 kT}{4\pi^2 r_0^2 \omega_0^2}. \quad (41.12)$$

Then we substitute the formula (41.6) for gamma (do not worry about writing ω_0 ; since it is true of any ω_0 , we may just call it ω) and the formula for $I(\omega)$ then comes out

$$I(\omega) = \frac{\omega^2 kT}{\pi^2 c^2}. \quad (41.13)$$

And that gives us the distribution of light in a hot furnace. It is called the *black-*

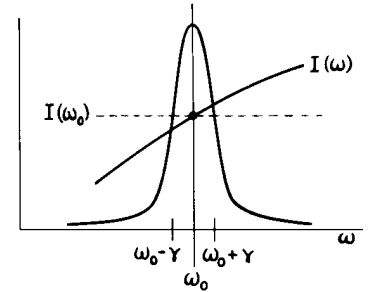


Fig. 41-3. The factors in the integrand (41.10). The peak is the resonance curve $1/(\omega - \omega_0)^2 + \gamma^2/4$. To a good approximation the factor $I(\omega)$ can be replaced by $I(\omega_0)$.

body radiation. Black, because the hole in the furnace that we look at is black when the temperature is zero.

Inside a closed box at temperature T , (41.13) is the distribution of energy of the radiation, according to classical theory. First, let us notice a remarkable feature of that expression. The charge of the oscillator, the mass of the oscillator, all properties specific to the oscillator, *cancel out*, because once we have reached equilibrium with one oscillator, we must be at equilibrium with any other oscillator of a different mass, or we will be in trouble. So this is an important kind of check on the proposition that equilibrium does not depend on what we are in equilibrium with, but *only on the temperature*. Now let us draw a picture of the $I(\omega)$ curve (Fig. 41-4). It tells us how much light we have at different frequencies.

The amount of intensity that there is in our box, per unit frequency range, goes, as we see, as the square of the frequency, which means that if we have a box at any temperature at all, and if we look at the x-rays that are coming out, there will be a lot of them!

Of course we know this is false. When we open the furnace and take a look at it, we do not burn our eyes out from x-rays at all. It is completely false. Furthermore, the *total energy* in the box, the total of all this intensity summed over all frequencies, would be the area under this infinite curve. Therefore, something is fundamentally, powerfully, and absolutely wrong.

Thus was the classical theory *absolutely incapable* of correctly describing the distribution of light from a blackbody, just as it was incapable of correctly describing the specific heats of gases. Physicists went back and forth over this derivation from many different points of view, and there is no escape. This is the prediction of classical physics. Equation (41.13) is called *Rayleigh's law*, and it is the prediction of classical physics, and is obviously absurd.

41-3 Equipartition and the quantum oscillator

The difficulty above was another part of the continual problem of classical physics, which started with the difficulty of the specific heat of gases, and now has been focused on the distribution of light in a blackbody. Now, of course, at the time that theoreticians studied this thing, there were also many *measurements* of the actual curve. And it turned out that the correct curve looked like the dashed curves in Fig. 41-4. That is, the x-rays were not there. If we lower the temperature, the whole curve goes down in proportion to T , according to the classical theory, but the observed curve also cuts off sooner at a lower temperature. Thus the low-frequency end of the curve is right, but the high-frequency end is wrong. Why? When Sir James Jeans was worrying about the specific heats of gases, he noted that motions which have high frequency are "frozen out" as the temperature goes too low. That is, if the temperature is too low, if the frequency is too high, the oscillators *do not have* kT of energy on the average. Now recall how our derivation of (41.13) worked: It all depends on the energy of an oscillator at thermal equilibrium. What the kT of (41.5) was, and what the same kT in (41.13) is, is the mean energy of a harmonic oscillator of frequency ω at temperature T . Classically, this is kT , but experimentally, no!—not when the temperature is too low or the oscillator frequency is too high. And so the reason that the curve falls off is the same reason that the specific heats of gases fail. It is easier to study the blackbody curve than it is the specific heats of gases, which are so complicated, therefore our attention is focused on determining the true blackbody curve, because this curve is a curve which correctly tells us, at every frequency, what the average energy of harmonic oscillators actually is as a function of temperature.

Planck studied this curve. He first determined the answer empirically, by fitting the observed curve with a nice function that fitted very well. Thus he had an empirical formula for the average energy of a harmonic oscillator as a function of frequency. In other words, he had the *right* formula instead of kT , and then by fiddling around he found a simple derivation for it which involved a very peculiar assumption. That assumption was that *the harmonic oscillator can take up energies only $\hbar\omega$ at a time*. The idea that they can have *any energy at all* is false. Of course, that was the beginning of the end of classical mechanics.

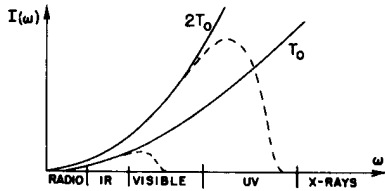


Fig. 41-4. The blackbody intensity distribution at two temperatures, according to classical physics (solid curves). The dashed curves show the actual distribution.

The very first correctly determined quantum-mechanical formula will now be derived. Suppose that the permitted energy levels of a harmonic oscillator were equally spaced at $\hbar\omega_0$ apart, so that the oscillator could take on only these different energies (Fig. 41-5). Planck made a somewhat more complicated argument than the one that is being given here, because that was the very beginning of quantum mechanics and he had to prove some things. But we are going to take it as a fact (which he demonstrated in this case) that the probability of occupying a level of energy E is $P(E) = \alpha e^{-E/kT}$. If we go along with that, we will obtain the right result.

Suppose now that we have a lot of oscillators, and each is a vibrator of frequency ω_0 . Some of these vibrators will be in the bottom quantum state, some will be in the next one, and so forth. What we would like to know is the average energy of all these oscillators. To find out, let us calculate the total energy of all the oscillators and divide by the number of oscillators. That will be the average energy per oscillator in thermal equilibrium, and will also be the energy that is in equilibrium with the blackbody radiation and that should go in Eq. (41.13) in place of kT . Thus we let N_0 be the number of oscillators that are in the ground state (the lowest energy state); N_1 the number of oscillators in the state E_1 ; N_2 the number that are in state E_2 ; and so on. According to the hypothesis (which we have not proved) that in quantum mechanics the law that replaced the probability $e^{-P.E./kT}$ or $e^{-K.E./kT}$ in classical mechanics is that the probability goes down as $e^{-\Delta E/kT}$, where ΔE is the excess energy, we shall assume that the number N_1 that are in the first state will be the number N_0 that are in the ground state, times $e^{-\hbar\omega/kT}$. Similarly, N_2 , the number of oscillators in the second state, is $N_2 = N_0 e^{-2\hbar\omega/kT}$. To simplify the algebra, let us call $e^{-\hbar\omega/kT} = x$. Then we simply have $N_1 = N_0 x$, $N_2 = N_0 x^2$, \dots , $N_n = N_0 x^n$.

The total energy of all the oscillators must first be worked out. If an oscillator is in the ground state, there is no energy. If it is in the first state, the energy is $\hbar\omega_0$, and there are N_1 of them. So $N_1 \hbar\omega$, or $\hbar\omega N_0 x$ is how much energy we get from those. Those that are in the second state have $2\hbar\omega_0$, and there are N_2 of them, so $N_2 \cdot 2\hbar\omega = 2\hbar\omega N_0 x^2$ is how much energy we get, and so on. Then we add it all together to get $E_{\text{tot}} = N_0 \hbar\omega (0 + x + 2x^2 + 3x^3 + \dots)$.

And now, how many oscillators are there? Of course, N_0 is the number that are in the ground state, N_1 in the first state, and so on, and we add them together: $N_{\text{tot}} = N_0 (1 + x + x^2 + x^3 + \dots)$. Thus the average energy is

$$\langle E \rangle = \frac{E_{\text{tot}}}{N_{\text{tot}}} = \frac{N_0 \hbar\omega (0 + x + 2x^2 + 3x^3 + \dots)}{N_0 (1 + x + x^2 + \dots)}. \quad (41.14)$$

Now the two sums which appear here we shall leave for the reader to play with and have some fun with. When we are all finished summing and substituting for x in the sum, we should get—if we make no mistakes in the sum—

$$\langle E \rangle = \frac{\hbar\omega}{e^{\hbar\omega/kT} - 1}. \quad (41.15)$$

This, then, was the first quantum-mechanical formula ever known, or ever discussed, and it was the beautiful culmination of decades of puzzlement. Maxwell knew that there was something wrong, and the problem was, what was *right*? Here is the quantitative answer of what is right instead of kT . This expression should, of course, approach kT as $\omega \rightarrow 0$ or as $T \rightarrow \infty$. See if you can prove that it does—learn how to do the mathematics.

This is the famous cutoff factor that Jeans was looking for, and if we use it instead of kT in (41.13), we obtain for the distribution of light in a black box

$$I(\omega) d\omega = \frac{\hbar\omega^3 d\omega}{\pi^2 c^2 (e^{\hbar\omega/kT} - 1)}. \quad (41.16)$$

We see that for a large ω , even though we have ω^3 in the numerator, there is an e raised to a tremendous power in the denominator, so the curve comes down again and does not “blow up”—we do not get ultraviolet light and x-rays where we do not expect them!

N_4	$E_4 = 4\hbar\omega$	$P_4 = A \exp(-4\hbar\omega/kT)$
N_3	$E_3 = 3\hbar\omega$	$P_3 = A \exp(-3\hbar\omega/kT)$
N_2	$E_2 = 2\hbar\omega$	$P_2 = A \exp(-2\hbar\omega/kT)$
N_1	$E_1 = \hbar\omega$	$P_1 = A \exp(-\hbar\omega/kT)$
N_0	$E_0 = 0$	$P_0 = A$

Fig. 41-5. The energy levels of a harmonic oscillator are equally spaced: $E_n = n\hbar\omega$.

One might complain that in our derivation of (41.16) we used the quantum theory for the energy levels of the harmonic oscillator, but the classical theory in determining the cross section σ_s . But the quantum theory of light interacting with a harmonic oscillator gives exactly the same result as that given by the classical theory. That, in fact, is why we were justified in spending so much time on our analysis of the index of refraction and the scattering of light, using a model of atoms like little oscillators—the quantum formulas are substantially the same.

Now let us return to the Johnson noise in a resistor. We have already remarked that the theory of this noise power is really the same theory as that of the classical blackbody distribution. In fact, rather amusingly, we have already said that if the resistance in a circuit were not a real resistance, but were an antenna (an antenna acts like a resistance because it radiates energy), a radiation resistance, it would be easy for us to calculate what the power would be. It would be just the power that runs into the antenna from the light that is all around, and we would get the same distribution, changed by only one or two factors. We can suppose that the resistor is a generator with an unknown power spectrum $P(\omega)$. The spectrum is determined by the fact that this same generator, connected to a resonant circuit of *any frequency*, as in Fig. 41-2(b), generates in the inductance a voltage of the magnitude given in Eq. (41.2). One is thus led to the same integral as in (41.10), and the same method works to give Eq. (41.3). For low temperatures the kT in (41.3) must of course be replaced by (41.15). The two theories (blackbody radiation and Johnson noise) are also closely related physically, for we may of course connect a resonant circuit to an *antenna*, so the resistance R is a pure *radiation resistance*. Since (41.2) does not depend on the physical origin of the resistance, we know the generator G for a real resistance and for radiation resistance is the same. What is the origin of the generated power $P(\omega)$ if the resistance R is only an ideal antenna in equilibrium with its environment at temperature T ? It is the radiation $I(\omega)$ in the space at temperature T which impinges on the antenna and, as “received signals,” makes an effective generator. Therefore one can deduce a direct relation of $P(\omega)$ and $I(\omega)$, leading then from (41.13) to (41.3).

All the things we have been talking about—the so-called Johnson noise and Planck’s distribution, and the correct theory of the Brownian movement which we are about to describe—are developments of the first decade or so of this century. Now with those points and that history in mind, we return to the Brownian movement.

41-4 The random walk

Let us consider how the position of a jiggling particle should change with time, for very long times compared with the time between “kicks.” Consider a little Brownian movement particle which is jiggling about because it is bombarded on all sides by irregularly jiggling water molecules. Query: After a given length of time, how far away is it likely to be from where it began? This problem was solved by Einstein and Smoluchowski. If we imagine that we divide the time into little intervals, let us say a hundredth of a second or so, then after the first hundredth of a second it moves here, and in the next hundredth it moves some more, in the next hundredth of a second it moves somewhere else, and so on. In terms of the rate of bombardment, a hundredth of a second is a very long time. The reader may easily verify that the number of collisions a single molecule of water receives in a second is about 10^{14} , so in a hundredth of a second it has 10^{12} collisions, which is a lot! Therefore, after a hundredth of a second it is not going to remember what happened before. In other words, the collisions are all *random*, so that one “step” is not related to the previous “step.” It is like the famous drunken sailor problem: the sailor comes out of the bar and takes a sequence of steps, but each step is chosen at an arbitrary angle, at random (Fig. 41-6). The question is: After a long time, where is the sailor? Of course we do not know! It is impossible to say. What do we mean—he is just somewhere more or less random. Well then, on the average, where is he? *On the average, how far away from the bar has he gone?* We have already answered this question, because once we were discussing the

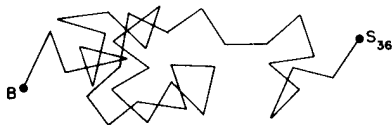


Fig. 41-6. A random walk of 36 steps of length l . How far is S_{36} from B ?
Ans: about $6l$ on the average.

superposition of light from a whole lot of different sources at different phases, and that meant adding a lot of arrows at different angles (Chapter 32). There we discovered that the mean square of the distance from one end to the other of the chain of random steps, which was the intensity of the light, is the sum of the intensities of the separate pieces. And so, by the same kind of mathematics, we can prove immediately that if \mathbf{R}_N is the vector distance from the origin after N steps, the mean square of the distance from the origin is proportional to the number N of steps. That is, $\langle R_N^2 \rangle = NL^2$, where L is the length of each step. Since the number of steps is proportional to the time in our present problem, *the mean square distance is proportional to the time*:

$$\langle R^2 \rangle = \alpha t. \quad (41.17)$$

This does not mean that the *mean distance* is proportional to the time. If the mean distance were proportional to the time it would mean that the drifting is at a nice uniform velocity. The sailor *is* making some relatively sensible headway, but only such that his mean *square* distance is proportional to time. That is the characteristic of a random walk.

We may show very easily that in each successive step the square of the distance increases, on the average, by L^2 . For if we write $\mathbf{R}_N = \mathbf{R}_{N-1} + \mathbf{L}$, we find that \mathbf{R}_N^2 is

$$\mathbf{R}_N \cdot \mathbf{R}_N = R_N^2 = R_{N-1}^2 + 2\mathbf{R}_{N-1} \cdot \mathbf{L} + L^2,$$

and averaging over many trials, we have $\langle R_N^2 \rangle = \langle R_{N-1}^2 \rangle + L^2$, since $\langle \mathbf{R}_{N-1} \cdot \mathbf{L} \rangle = 0$. Thus, by induction,

$$R_N^2 = NL^2. \quad (41.18)$$

Now we would like to calculate the coefficient α in Eq. (41.17), and to do so we must add a feature. We are going to suppose that if we were to put a force on this particle (having nothing to do with the Brownian movement—we are taking a side issue for the moment), then it would react in the following way against the force. First, there would be inertia. Let m be the coefficient of inertia, the effective mass of the object (not necessarily the same as the real mass of the real particle, because the water has to move around the particle if we pull on it). Thus if we talk about motion in one direction, there is a term like $m(d^2x/dt^2)$ on one side. And next, we want also to assume that if we kept a steady pull on the object, there would be a drag on it from the fluid, proportional to its velocity. Besides the inertia of the fluid, there is a resistance to flow due to the viscosity and the complexity of the fluid. It is absolutely essential that there *be* some irreversible losses, something like resistance, in order that there be fluctuations. There is no way to produce the kT unless there are also losses. The source of the fluctuations is very closely related to these losses. What the mechanism of this drag is, we will discuss soon—we shall talk about forces that are proportional to the velocity and where they come from. But let us suppose for now that there is such a resistance. Then the formula for the motion under an external force, when we are pulling on it in a normal manner, is

$$m \frac{d^2x}{dt^2} + \mu \frac{dx}{dt} = F_{\text{ext}}. \quad (41.19)$$

The quantity μ can be determined directly from experiment. For example, we can watch the drop fall under gravity. Then we know that the force is mg , and μ is mg divided by the speed of fall the drop ultimately acquires. Or we could put the drop in a centrifuge and see how fast it sediments. Or if it is charged, we can put an electric field on it. So μ is a measurable thing, not an artificial thing, and it is known for many types of colloidal particles, etc.

Now let us use the same formula in the case where the force is not external, but is equal to the irregular forces of the Brownian movement. We shall then try to determine the mean square distance that the object goes. Instead of taking the distances in three dimensions, let us take just one dimension, and find the mean of x^2 , just to prepare ourselves. (Obviously the mean of x^2 is the same as the mean of y^2 is the same as the mean of z^2 , and therefore the mean square of the distance

is just 3 times what we are going to calculate.) The x -component of the irregular forces is, of course, just as irregular as any other component. What is the rate of change of x^2 ? It is $d(x^2)/dt = 2x(dx/dt)$, so what we have to find is the average of the position times the velocity. We shall show that this is a constant, and that therefore the mean square radius will increase proportionally to the time, and at what rate. Now if we multiply Eq. (41.19) by x , $mx(d^2x/dt^2) + \mu x(dx/dt) = xF_x$. We want the time average of $x(dx/dt)$, so let us take the average of the whole equation, and study the three terms. Now what about x times the force? If the particle happens to have gone a certain distance x , then, since the irregular force is *completely* irregular and does not know where the particle started from, the next impulse can be in any direction relative to x . If x is positive, there is no reason why the average force should also be in that direction. It is just as likely to be one way as the other. The bombardment forces are not driving it in a definite direction. So the average value of x times F is zero. On the other hand, for the term $mx(d^2x/dt^2)$ we will have to be a little fancy, and write this as

$$mx \frac{d^2x}{dt^2} = m \frac{d[x(dx/dt)]}{dt} - m \left(\frac{dx}{dt} \right)^2.$$

Thus we put in these two terms and take the average of both. So let us see how much x times the velocity should be. Now x times the velocity has a mean that does not change with time, because when it gets to some position it has no remembrance of where it was before, so things are no longer changing with time. So this quantity, on the average, is zero. We have left the quantity mv^2 , and that is the only thing we know: $mv^2/2$ has a mean value $\frac{1}{2}kT$. Therefore we find that

$$\left\langle mx \frac{d^2x}{dt^2} \right\rangle + \mu \left\langle x \frac{dx}{dt} \right\rangle = \langle xF_x \rangle$$

implies

$$-\langle mv^2 \rangle + \frac{\mu}{2} \frac{d}{dt} \langle x^2 \rangle = 0,$$

or

$$\frac{d\langle x^2 \rangle}{dt} = 2 \frac{kT}{\mu}. \quad (41.20)$$

Therefore the object has a mean square distance $\langle R^2 \rangle$, at the end of a certain amount of t , equal to

$$\langle R^2 \rangle = 6kT \frac{t}{\mu}. \quad (41.21)$$

And so we can actually determine *how far* the particles go! We first must determine how they react to a steady force, how fast they drift under a known force (to find μ), and then we can determine how far they go in their random motions. This equation was of considerable importance historically, because it was one of the first ways by which the constant k was determined. After all, we can measure μ , the time, how far the particles go, and we can take an average. The reason that the determination of k was important is that in the law $PV = RT$ for a mole, we know that R , which can also be measured, is equal to the number of atoms in a mole times k . A mole was originally defined as so and so many *grams* of oxygen—16 (now carbon is used), so the number of *atoms* in a mole was not known, originally. It is, of course, a very interesting and important problem. How big are atoms? How many are there? So one of the earliest determinations of the number of atoms was by the determination of how far a dirty little particle would move if we watched it patiently under a microscope for a certain length of time. And thus Boltzmann's constant k and the Avogadro number N_0 were determined because R had already been measured.

Applications of Kinetic Theory

42-1 Evaporation

In this chapter we shall discuss some further applications of kinetic theory. In the previous chapter we emphasized one particular aspect of kinetic theory, namely, that the average kinetic energy in any degree of freedom of a molecule or other object is $\frac{1}{2}kT$. The central feature of what we shall now discuss, on the other hand, is the fact that the probability of finding a particle in different places, per unit volume, varies as $e^{-\text{potential energy}/kT}$; we shall make a number of applications of this.

The phenomena which we want to study are relatively complicated: a liquid evaporating, or electrons in a metal coming out of the surface, or a chemical reaction in which there are a large number of atoms involved. In such cases it is no longer possible to make from the kinetic theory any simple and correct statements, because the situation is too complicated. Therefore, this chapter, except where otherwise emphasized, is quite inexact. The idea to be emphasized is only that we can understand, from the kinetic theory, *more or less* how things ought to behave. By using thermodynamic arguments, or some empirical measurements of certain critical quantities, we can get a more accurate representation of the phenomena.

However, it is very useful to know even only more or less why something behaves as it does, so that when the situation is a new one, or one that we have not yet started to analyze, we can say, more or less, what ought to happen. So this discussion is highly inaccurate but essentially right—right in idea, but a little bit simplified, let us say, in the specific details.

The first example that we shall consider is the evaporation of a liquid. Suppose we have a box with a large volume, partially filled with liquid in equilibrium and with the vapor at a certain temperature. We shall suppose that the molecules of the vapor are relatively far apart, and that inside the liquid, the molecules are packed close together. The problem is to find out how many molecules there are in the vapor phase, compared with the number there are in the liquid. How dense is the vapor at a given temperature, and how does it depend on the temperature?

Let us say that n equals the number of molecules per unit volume in the vapor. That number, of course, varies with the temperature. If we add heat, we get more evaporation. Now let another quantity, $1/V_a$, equal the number of atoms per unit volume in the liquid: We suppose that each molecule in the liquid occupies a certain volume, so that if there are more molecules of liquid, then all together they occupy a bigger volume. Thus if V_a is the volume occupied by one molecule, the number of molecules in a unit volume is a unit volume divided by the volume of each molecule. Furthermore, we suppose that there is a force of attraction between the molecules to hold them together in the liquid. Otherwise we cannot understand why it condenses. Thus suppose that there is such a force and that there is an energy of binding of the molecules in the liquid which is lost when they go into the vapor. That is, we are going to suppose that, in order to take a single molecule out of the liquid into the vapor, a certain amount of work W has to be done. There is a certain difference, W , in the energy of a molecule in the liquid from what it would have if it were in the vapor, because we have to pull it away from the other molecules which attract it.

Now we use the general principle that the number of atoms per unit volume in two different regions is $n_2/n_1 = e^{-(E_2-E_1)/kT}$. So the number n per unit volume in the vapor, divided by the number $1/V_a$ per unit volume in the liquid, is equal to

$$nV_a = e^{-W/kT}, \quad (42.1)$$

42-1 Evaporation

42-2 Thermionic emission

42-3 Thermal ionization

42-4 Chemical kinetics

42-5 Einstein's laws of radiation

because that is the general rule. It is like the atmosphere in equilibrium under gravity, where the gas at the bottom is denser than that at the top because of the work mgh needed to lift the gas molecules to the height h . In the liquid, the molecules are denser than in the vapor because we have to pull them out through the energy "hill" W , and the ratio of the densities is $e^{-W/kT}$.

This is what we wanted to deduce—that the vapor density varies as e to the minus some energy or other over kT . The factors in front are not really interesting to us, because in most cases the vapor density is very much lower than the liquid density. In those circumstances, where we are not near the critical point where they are almost the same, but where the vapor density is much lower than the liquid density, then the fact that n is very much less than $1/V_a$ is occasioned by the fact that W is very much greater than kT . So formulas such as (42.1) are interesting only when W is very much bigger than kT , because in those circumstances, since we are raising e to minus a tremendous amount, if we change T a little bit, that tremendous power changes a bit, and the change produced in the exponential factor is very much more important than any change that might occur in the factors out in front. Why should there be any changes in such factors as V_a ? Because ours was an approximate analysis. After all, there is not really a definite volume for each molecule; as we change the temperature, the volume V_a does not stay constant—the liquid expands. There are other little features like that, and so the actual situation is more complicated. There are slowly varying temperature-dependent factors all over the place. In fact, we might say that W itself varies slightly with temperature, because at a higher temperature, at a different molecular volume, there would be different average attractions, and so on. So, while we might think that if we have a formula in which everything varies in an unknown way with temperature then we have no formula at all, if we realize that the exponent W/kT is, in general, very large, we see that in the curve of the vapor density as a function of temperature most of the variation is occasioned by the exponential factor, and if we take W as a constant and the coefficient $1/V_a$ as nearly constant, it is a good approximation for short intervals along the curve. Most of the variation, in other words, is of the general nature $e^{-W/kT}$.

It turns out that there are many, many phenomena in nature which are characterized by having to borrow an energy from somewhere, and in which the central feature of the temperature variation is e to the minus the energy over kT . This is a useful fact only when the energy is large compared with kT , so that most of the variation is contained in the variation of the kT and not in the constant and in other factors.

Now let us consider another way of obtaining a somewhat similar result for the evaporation, but looking at it in more detail. To arrive at (42.1), we simply applied a rule which is valid at equilibrium, but in order to understand things better, there is no harm in trying to look at the details of what is going on. We may also describe what is going on in the following way: the molecules that are in the vapor continually bombard the surface of the liquid; when they hit it, they may bounce off or they may get stuck. There is an unknown factor for that—maybe 50–50, maybe 10 to 90—we do not know. Let us say they always get stuck—we can analyze it over again later on the assumption that they do not always get stuck. Then at a given moment there will be a certain number of atoms which are condensing onto the surface of the liquid. The number of condensing molecules, the number that arrive on a unit area, is the number n per unit volume times the velocity v . This velocity of the molecules is related to the temperature, because we know that $\frac{1}{2}mv^2$ is equal to $\frac{3}{2}kT$ on the average. So v is some kind of a mean velocity. Of course we should integrate over the angles and get some kind of an average, but it is roughly proportional to the root-mean-square velocity, within some factor. Thus

$$N_c = nv \quad (42.2)$$

is the number which arrive per unit area and are condensing.

At the same time, however, the atoms in the liquid are jiggling about, and from time to time one of them gets kicked out. Now we have to estimate how fast they get kicked out. The idea will be that at equilibrium the number that are kicked out per second and the number that arrive per second are equal.

How many get kicked out? In order to get kicked out, a particular molecule has to have acquired by accident an excess energy over its neighbors—a considerable excess energy, because it is attracted very strongly by the other molecules in the liquid. Ordinarily it does not leave because it is so strongly attracted, but in the collisions sometimes one of them gets an extra energy by accident. And the chance that it gets the extra energy W which it needs in our case is very small if $W \gg kT$. In fact, $e^{-W/kT}$ is the chance that an atom has picked up more than this much energy. That is the general principle in kinetic theory: in order to borrow an excess energy W over the average, the odds are e to the minus the energy that we have to borrow, over kT . Now suppose that some molecules have borrowed this energy. We now have to estimate how many leave the surface per second. Of course, just because a molecule has the necessary energy does not mean that it will actually evaporate, since it may be buried too deeply inside the liquid or, even if it is near the surface, it may be travelling in the wrong direction. The number that are going to leave a unit area per second is going to be something like this: the number of atoms there are near the surface, per unit area, divided by the time it takes one to escape, multiplied by the probability $e^{-W/kT}$ that they are ready to escape in the sense that they have enough energy.

We shall suppose that each molecule at the surface of the liquid occupies a certain cross-sectional area A . Then the number of molecules per unit area of liquid surface will be $1/A$. And now, how long does it take a molecule to escape? If the molecules have a certain average speed v , and have to move, say, one molecular diameter D , the thickness of the first layer, then the time it takes to get across that thickness is the time needed to escape, if the molecule has enough energy. The time will be D/v . Thus the number evaporating should be approximately

$$N_e = (1/A)(v/D)e^{-W/kT}. \quad (42.3)$$

Now the area of each atom times the thickness of the layer is approximately the same as the volume V_a occupied by a single atom. And so, in order to get equilibrium, we must have $N_e = N_c$, or

$$nv = (v/V_a)e^{-W/kT}. \quad (42.4)$$

We may cancel the v 's, since they are equal; even though one is the velocity of a molecule in the vapor and the other is the velocity of an evaporating molecule, these are the same, because we know their mean kinetic energy (in one direction) is $\frac{1}{2}kT$. But one may object, "No! No! These are the especially fast-moving ones; these are the ones that have picked up excess energy." Not really, because the moment they start to pull away from the liquid, they have to *lose* that excess energy against the potential energy. So, as they come to the surface they are slowed down to the velocity v ! It is the same as it was in our discussion of the distribution of molecular velocities in the atmosphere—at the bottom, the molecules had a certain distribution of energy. The ones that arrive at the top have the *same* distribution of energy, because the slow ones did not arrive at all, and the fast ones were slowed down. The molecules that are evaporating have the same distribution of energy as the ones inside—a rather remarkable fact. Anyway, it is useless to try to argue so closely about our formula because of other inaccuracies, such as the probability of bouncing back rather than entering the liquid, and so on. Thus we have a rough idea of the rate of evaporation and condensation, and we see, of course, that the vapor density n varies in the same way as before, but now we have understood it in some detail rather than just as an arbitrary formula.

This deeper understanding permits us to analyze some things. For example, suppose that we were to pump away the vapor at such a great rate that we removed the vapor as fast as it formed (if we had very good pumps and the liquid was evaporating very slowly), how fast would evaporation occur if we maintained a liquid temperature T ? Suppose that we have already experimentally measured the equilibrium vapor density, so that we know, at the given temperature, how many molecules per unit volume are in equilibrium with the liquid. Now we would like to know *how fast* it will evaporate. Even though we have used only a rough analysis so far as the evaporation part of it is concerned, the number of vapor

molecules *arriving* was not done so badly, aside from the unknown factor of reflection coefficient. So therefore we may use the fact that the number that are leaving, at equilibrium, is the same as the number that arrive. True, the vapor is being swept away and so the molecules are only coming out, but if the vapor were left alone, it would attain the equilibrium density at which the number that come back would equal the number that are evaporating. Therefore, we can easily see that the number that are coming off the surface per second is equal to the unknown reflection coefficient R times the number that would come down to the surface per second were the vapor still there, because that is how many would balance the evaporation at equilibrium:

$$N_e = nvR = (vR/V_a)e^{-W/kT}. \quad (42.5)$$

Of course, the number of molecules that hit the liquid from the vapor is easy to calculate, since we do not need to know as much about the forces as we do when we are worrying about how they get to escape through the liquid surface; it is much easier to make the argument the other way.

42-2 Thermionic emission

We may give another example of a very practical situation that is similar to the evaporation of a liquid—so similar that it is not worth making a separate analysis. It is essentially the same problem. In a radio tube there is a source of electrons, namely a heated tungsten filament, and a positively charged plate to attract the electrons. Any electron that escapes from the surface of the tungsten is immediately swept away to the plate. That is our ideal “pump,” which is “pumping” the electrons away all the time. Now the question is: How many electrons per second can we get out of a piece of tungsten, and how does that number vary with temperature? The answer to that problem is the same as (42.5), because it turns out that in a piece of metal, electrons are attracted to the ions, or to atoms, of the metal. They are attracted, if we may say it crudely, to the metal. In order to get an electron out of a piece of metal, it takes a certain amount of energy or work to pull it out. This work varies with the different kinds of metal. In fact, it varies even with the character of the surface of a given kind of metal, but the total work may be a few electron volts, which, incidentally, is typical of the energy involved in chemical reactions. We can remember the latter fact by remembering that the voltage in a chemical cell like a flashlight battery, which is produced by chemical reactions, is about one volt.

How can we find out how many electrons come out per second? It would be quite difficult to analyze the effects on the electrons going out; it is easier to analyze the situation the other way. So, we could start out by imagining that we did not draw the electrons away, and that the electrons were like a gas, and could come back to the metal. Then there would be a certain density of electrons at equilibrium which would, of course, be given by exactly the same formula as (42.1), where V_a is the volume per electron in the metal, roughly, and W is equal to $q_e\phi$, where ϕ is the so-called *work function*, or the voltage needed to pull an electron off the surface. This would tell us how many electrons would have to be in the surrounding space and striking the metal in order to balance the ones that are coming out. And thus it is easy to calculate how many are coming out if we sweep away all of them, because the number that are coming out is exactly equal to the number that would be going in with the above density of electron “vapor.” In other words, the answer is that the current of electricity that comes in per unit area is equal to the charge on each times the number that arrive per second per unit area, which is the number per unit volume times the velocity, as we have seen many times:

$$I = q_e nv = (q_e v R / V_a) e^{-q_e \phi / kT}. \quad (42.6)$$

Now one electron volt corresponds to kT at a temperature of 11,600 degrees. The filament of the tube may be operating at a temperature of, say, 1100 degrees, so the exponential factor is something like e^{-10} ; when we change the temperature a

little bit, the exponential factor changes a lot. Thus, again, the central feature of the formula is the $e^{-q_e\phi/kT}$. As a matter of fact, the factor in front is quite wrong—it turns out that the behavior of electrons in a metal is not correctly described by the classical theory, but by quantum mechanics, but this only changes the factor in front a little. Actually, no one has ever been able to get the thing straightened out very well, even though many people have used the high-class quantum-mechanical theory for their calculations. The big problem is, does W change slightly with temperature? If it does, one cannot distinguish a W changing slowly with temperature from a different coefficient in front. That is, if W changed linearly, say, with temperature, so that $W = W_0 + \alpha kT$, then we would have

$$e^{-W/kT} = e^{-(W_0 + \alpha kT)/kT} = e^{-\alpha} e^{-W_0/kT}.$$

Thus a linearly temperature-dependent W is equivalent to a shifted “constant.” It is really quite difficult and usually fruitless to try to obtain the coefficient in the front accurately.

42-3 Thermal ionization

Now we go on to another example of the same idea; always the same idea. This has to do with ionization. Suppose that in a gas we have a whole lot of atoms which are in the neutral state, say, but the gas is hot and the atoms can become ionized. We would like to know how many ions there are in a given circumstance if we have a certain density of atoms per unit volume at a certain temperature. Again we consider a box in which there are N atoms which can hold electrons. (If an electron has come off an atom, it is called an *ion*, and if the atom is neutral, we simply call it an atom.) Then suppose that, at any given moment, the number of neutral atoms is n_a , the number of ions is n_i , and the number of electrons is n_e , all per unit volume. The problem is: What is the relationship of these three numbers?

In the first place, we have two conditions or constraints on the numbers. For instance, as we vary different conditions, like the temperature and so on, $n_a + n_i$ would remain constant, because this would be simply the number N of atomic nuclei that are in the box. If we keep the number of nuclei per unit volume fixed, and change, say, the temperature, then as the ionization proceeded some atoms would turn to ions, but the total number of atoms plus ions would be unchanged. That is, $n_a + n_i = N$. Another condition is that if the entire gas is to be electrically neutral (and if we neglect double or triple ionization), that means that the number of ions is equal to the number of electrons at all times, or $n_i = n_e$. These are subsidiary equations that simply express the conservation of charge and the conservation of atoms.

These equations are true, and we ultimately will use them when we consider a real problem. But we want to obtain another relationship between the quantities. We can do this as follows. We again use the idea that it takes a certain amount of energy to lift the electron out of the atom, which we call the *ionization energy*, and we would write it as W , in order to make all of the formulas look the same. So we let W equal the energy needed to pull an electron out of an atom and make an ion. Now we again say that the number of free electrons per unit volume in the “vapor” is equal to the number of bound electrons per unit volume in the atoms, times e to the minus the energy difference between being bound and being free, over kT . That is the basic equation again. How can we write it? The number of free electrons per unit volume would, of course, be n_e , because that is the definition of n_e . Now what about the number of electrons per unit volume that are bound to atoms? The total number of places that we could put the electrons is apparently $n_a + n_i$, and we will suppose that when they are bound each one is bound within a certain volume V_a . So the total amount of volume which is available to electrons which would be bound is $(n_a + n_i)V_a$, so we might want to write our formula as

$$n_e = \frac{n_a}{(n_a + n_i)V_a} e^{-W/kT}.$$

The formula is wrong, however, in one essential feature, which is the following: when an electron is already on an atom, another electron cannot come to that volume anymore! In other words, all the volumes of all the possible sites are not really available for the one electron which is trying to make up its mind whether or not to be in the vapor or in the condensed position, because in this problem there is an extra feature that when one electron is where another electron is, it is not allowed to go—it is repelled. For that reason, it comes out that we should count only that part of the volume which is available for an electron to sit on or not. That is, those which are already occupied do not count in the total available volume, but the only volume which is allowed is that of the *ions*, where there are vacant places for the electron to go. Then, in those circumstances, we find that a nicer way to write our formula is

$$\frac{n_e n_i}{n_a} = \frac{1}{V_a} e^{-W/kT}. \quad (42.7)$$

This formula is called the *Saha ionization equation*. Now let us see if we can understand qualitatively why a formula like this is right, by arguing about the kinetic things that are happening.

First, every once in a while an electron comes to an ion and they combine to make an atom. And also, every once in a while, an atom gets into a collision and breaks up into an ion and an electron. Now those two rates must be equal. How fast do electrons and ions find each other? The rate is certainly increased if the number of electrons per unit volume is increased. It is also increased if the number of ions per unit volume is increased. That is, the total rate at which recombination is occurring is certainly proportional to the number of electrons times the number of ions. Now the total rate at which ionization is occurring due to collisions must be dependent linearly on how many atoms there are to ionize. And so the rates will balance when there is some relationship between the product $n_e n_i$ and the number of atoms, n_a . The fact that this relationship happens to be given by this particular formula, where W is the ionization energy, is of course a little bit more information, but we can easily understand that the formula would necessarily involve the concentrations of the electrons, ions, and atoms in the combination $n_e n_i / n_a$ to produce a constant independent of the n 's, and dependent only on temperature, the atomic cross sections, and other constant factors.

We may also note that, since the equation involves the numbers *per unit volume*, if we were to do two experiments with a given total number N of atoms plus ions, that is, a certain fixed number of nuclei, but using boxes with different volumes, the n 's would all be smaller in the larger box. But since the ratio $n_e n_i / n_a$ stays the same, the *total number* of electrons and ions must be greater in the larger box. To see this, suppose that there are N nuclei inside a box of volume V , and that a fraction f of them are ionized. Then $n_e = fN/V = n_i$, and $n_a = (1-f)N/V$. Then our equation becomes

$$\frac{f^2}{1-f} \frac{N}{V} = \frac{e^{-W/kT}}{V_a}. \quad (42.8)$$

In other words, if we take a smaller and smaller density of atoms, or make the volume of the container bigger and bigger, the fraction f of electrons and ions must increase. That ionization, just from “expansion” as the density goes down, is the reason why we believe that at very low densities, such as in the cold space between the stars, there may be ions present, even though we might not understand it from the point of view of the available energy. Although it takes many, many kT of energy to make them, there are ions present.

Why can there be ions present when there is so much space around, while if we increase the density, the ions tend to disappear? *Answer*: Consider an atom. Every once in a while, light, or another atom, or an ion, or whatever it is that maintains thermal equilibrium, strikes it. Very rarely, because it takes such a terrific amount of excess energy, an electron comes off and an ion is left. Now that electron, if the space is enormous, wanders and wanders and does not come near anything for years, perhaps. But once in a very great while, it does come back

to an ion and they combine to make an atom. So the rate at which electrons are coming out from the atoms is very slow. But if the volume is enormous, an electron which has escaped takes so long to find another ion to recombine with that its probability of recombination is very, very small; thus, in spite of the large excess energy needed, there may be a reasonable number of electrons.

42-4 Chemical kinetics

The same situation that we have just called “ionization” is also found in a chemical reaction. For instance, if two objects A and B combine into a compound AB , then if we think about it for a while we see that AB is what we have called an atom, B is what we call an electron, and A is what we call an ion. With these substitutions the equations of equilibrium are exactly the same in form:

$$\frac{n_A n_B}{n_{AB}} = c e^{-W/kT}. \quad (42.9)$$

This formula, of course, is not exact, since the “constant” c depends on how much volume is allowed for the A and B to combine, and so on, but by thermodynamic arguments one can identify what the meaning of the W in the exponential factor is, and it turns out that it is very close to the energy needed in the reaction.

Suppose that we tried to understand this formula as a result of collisions, much in the way that we understood the evaporation formula, by arguing about how many electrons came off and how many of them came back per unit time. Suppose that A and B combine in a collision every once in a while to form a compound AB . And suppose that the compound AB is a complicated molecule which jiggles around and is hit by other molecules, and from time to time it gets enough energy to explode and break up again into A and B .

Now it actually turns out, in chemical reactions, that if the atoms come together with too small an energy, even though energy may be released in the reaction $A + B \rightarrow AB$, the fact that A and B may touch each other does not necessarily make the reaction start. It usually is required that the collision be rather hard, in fact, to get the reaction to go at all—a “soft” collision between A and B may not do it, even though energy may be released in the process. So let us suppose that it is very common in chemical reactions that, in order for A and B to form AB , they cannot just hit each other, but they have to hit each other *with sufficient energy*. This energy is called the *activation energy*—the energy needed to “activate” the reaction. Call A^* the activation energy, the excess energy needed in a collision in order that the reaction may really occur. Then the rate R_f at which A and B produce AB would involve the number of atoms of A times the number of atoms of B , times the rate at which a single atom would strike a certain cross section σ_{AB} , times a factor $e^{-A^*/kT}$, which is the probability that they have enough energy:

$$R_f = n_A n_B v \sigma_{AB} e^{-A^*/kT}. \quad (42.10)$$

Now we have to find the opposite rate, R_r . There is a certain chance that AB will fly apart. In order to fly apart, it not only must have the energy W which it needs in order to get apart at all but, just as it was hard for A and B to come together, so there is a kind of hill that A and B have to climb over to get apart again; they must have not only enough energy just to get ready to pull apart, but a certain excess. It is like climbing a hill to get into a deep valley; they have to climb the hill coming in and they have to climb out of the valley and then over the hill coming back (Fig. 42-1). Thus the rate at which AB goes to A and B will be proportional to the number n_{AB} that are present, times $e^{-(W+A^*)/kT}$:

$$R_r = c' n_{AB} e^{-(W+A^*)/kT}. \quad (42.11)$$

The c' will involve the volume of atoms and the rate of collisions, which we can work out, as we did the case of evaporation, with areas and times and thicknesses; but we shall not do this. The main feature of interest to us is that when these two

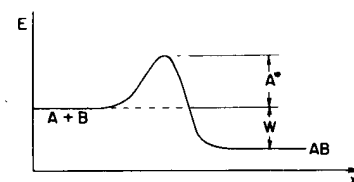


Fig. 42-1. The energy relationship for the reaction $A + B \rightarrow AB$.

rates are equal, the ratio of them is equal to unity. This tells us that $n_A n_B / n_{AB} = ce^{-W/kT}$, as before, where c involves the cross sections, velocities, and other factors independent of the n 's.

The interesting thing is that the rate of the reaction also varies as $e^{-\text{const}/kT}$, although the constant is not the same as that which governs the concentrations; the activation energy A^* is quite different from the energy W . W governs the proportions of A , B , and AB that we have in equilibrium, but if we want to know how fast $A + B$ goes to AB , that is not a question of equilibrium, and here a different energy, the activation energy, governs the rate of reaction through an exponential factor.

Furthermore, A^* is not a fundamental constant like W . Suppose that at the surface of the wall—or at some other place— A and B could temporarily stick there in such a way that they could combine more easily. In other words, we might find a “tunnel” through the hill, or perhaps a lower hill. By the conservation of energy, when we are all finished we have still made AB out of A and B , so the energy difference W will be quite independent of the way the reaction occurred, but the activation energy A^* will depend *very much* on the way the reaction occurs. This is why the rates of chemical reactions are very sensitive to outside conditions. We can change the rate by putting in a surface of a different kind, we can put it in a “different barrel” and it will go at a different rate, if it depends on the nature of the surface. Or if we put in a third kind of object it may change the rate very much; some things produce enormous changes in rate simply by changing the A^* a little bit—they are called *catalysts*. A reaction might practically not occur at all because A^* is too big at the given temperature, but when we put in this special stuff, the catalyst, then the reaction goes very fast indeed, because A^* is reduced.

Incidentally, there is some trouble with such a reaction, A plus B , making AB , because we cannot conserve both energy and momentum when we try to put two objects together to make one that is more stable. Therefore, we need at least a third object C , so the actual reaction is much more complicated. The forward rate would involve the product $n_A n_B n_C$, and it might seem that our formula is going wrong, but no! When we look at the rate at which AB goes the other way, we find that it also needs to collide with C , so there is an $n_{AB} n_C$ in the reverse rate; the n_C 's cancel out in the formula for the equilibrium concentrations. The law of equilibrium, (42.9), which we first wrote down is absolutely guaranteed to be true, no matter what the mechanism of the reaction may be!

42-5 Einstein's laws of radiation

We now turn to an interesting analogous situation having to do with the black-body radiation law. In the last chapter we worked out the distribution law for the radiation in a cavity the way Planck did, considering the radiation from an oscillator. The oscillator had to have a certain mean energy, and since it was oscillating, it would radiate and would keep pumping radiation into the cavity until it piled up enough radiation to balance the absorption and emission. In that way we found that the intensity of radiation at frequency ω was given by the formula

$$I(\omega) d\omega = \frac{\hbar \omega^3 d\omega}{\pi^2 c^2 (e^{\hbar \omega / kT} - 1)} . \quad (42.12)$$

This result involved the assumption that the oscillator which was generating the radiation had definite, equally spaced energy levels. We did not say that light had to be a photon or anything like that. There was no discussion about how, when an atom goes from one level to another, the energy must come out in one unit of energy, $\hbar \omega$, in the form of light. Planck's original idea was that the matter was quantized but not the light: material oscillators cannot take up just any energy, but have to take it in lumps. Furthermore, the trouble with the derivation is that it was partially classical. We calculated the rate of radiation from an oscillator according to classical physics; then we turned around and said, “No, this oscillator has a lot of energy levels.” So gradually, in order to find the right result, the completely quantum-mechanical result, there was a slow develop-

ment which culminated in the quantum mechanics of 1927. But in the meantime, there was an attempt by Einstein to convert Planck's viewpoint that only oscillators of matter were quantized, to the idea that light was really photons and could be considered in a certain way as particles with energy $\hbar\omega$. Furthermore, Bohr had pointed out that *any* system of atoms has energy levels, but they are not necessarily equally spaced like Planck's oscillator. And so it became necessary to rederive or at least rediscuss the radiation law from a more completely quantum-mechanical viewpoint.

Einstein assumed that Planck's final formula was right, and he used that formula to obtain some new information, previously unknown, about the interaction of radiation with matter. His discussion went as follows: Consider any two of the many energy levels of an atom, say the m th level and the n th level (Fig. 42-2). Now Einstein proposed that when such an atom has light of the right frequency shining on it, it can absorb that photon of light and make a transition from state n to state m , and that the probability that this occurs per second depends upon the two levels, of course, but is *proportional to how intense the light is* that is shining on it. Let us call the proportionality constant B_{nm} , merely to remind us that this is not a universal constant of nature, but depends on the particular pair of levels: some levels are easy to excite; some levels are hard to excite. Now what is the formula going to be for the rate of emission from m to n ? Einstein proposed that this must have two parts to it. First, even if there were no light present, there would be some chance that an atom in an excited state would fall to a lower state, emitting a photon; this we call *spontaneous emission*. It is analogous to the idea that an oscillator with a certain amount of energy, even in classical physics, does not keep that energy, but loses it by radiation. Thus the analog of spontaneous radiation of a classical system is that if the atom is in an excited state there is a certain probability A_{mn} , which depends on the levels again, for it to go down from m to n , and this probability is independent of whether light is shining on the atom or not. But then Einstein went further, and by comparison with the classical theory and by other arguments, concluded that emission was also influenced by the presence of light—that when light of the right frequency is shining on an atom, it has an increased rate of emitting a photon that is proportional to the intensity of the light, with a proportionality constant B_{mn} . Later, if we deduce that this coefficient is zero, then we will have found that Einstein was wrong. Of course we will find he was right.

Thus Einstein assumed that there are three kinds of processes: an absorption proportional to the intensity of light, an emission proportional to the intensity of light, called *induced emission* or sometimes *stimulated emission*, and a spontaneous emission independent of light.

Now suppose that we have, in equilibrium at temperature T , a certain number of atoms N_n in the state n and another number N_m in the state m . Then the total number of atoms that are going from n to m is the number that are in the state n times the rate per second that, if one is in n , it goes up to m . So we have a formula for the number that are going from n to m per second:

$$R_{n \rightarrow m} = N_n B_{nm} I(\omega). \quad (42.13)$$

The number that will go from m to n is expressed in the same manner, as the number N_m that are in m , times the chance per second that each one goes down to n . This time our expression is

$$R_{m \rightarrow n} = N_m [A_{mn} + B_{mn} I(\omega)]. \quad (42.14)$$

Now we shall suppose that in thermal equilibrium the number of atoms going up must equal the number coming down. That is one way, at least, in which the number will be sure to stay constant in each level.* So we take these two rates

* This is not the only way one can arrange to keep the numbers of atoms in the various levels constant, but it is the way it actually works. That every process must, in thermal equilibrium, be balanced by its exact opposite is called the *principle of detailed balancing*.

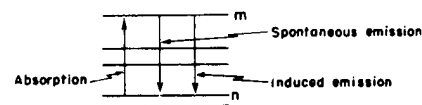


Fig. 42-2. Transitions between two energy levels of an atom.

to be equal at equilibrium. But we have one other piece of information: we know how large N_m is compared with N_n —the ratio of those two is $e^{-(E_m - E_n)/kT}$. Now Einstein assumed that the only light which is effective in making the transition from n to m is the light which has the frequency corresponding to the energy difference, so $E_m - E_n = \hbar\omega$ in all our formulas. Thus

$$N_m = N_n e^{-\hbar\omega/kT}. \quad (42.15)$$

Thus if we set the two rates equal: $N_n B_{nm} I(\omega) = N_m [A_{mn} + B_{mn} I(\omega)]$, and divide by N_m , we get

$$B_{nm} I(\omega) e^{\hbar\omega/kT} = A_{mn} + B_{mn} I(\omega). \quad (42.16)$$

From this equation, we can calculate $I(\omega)$. It is simply

$$I(\omega) = \frac{A_{mn}}{B_{nm} e^{\hbar\omega/kT} - B_{mn}}. \quad (42.17)$$

But Planck has already told us that the formula must be (42.12). Therefore we can deduce something: First, that B_{nm} must equal B_{mn} , since otherwise we cannot get the $(e^{\hbar\omega/kT} - 1)$. So Einstein discovered some things that he did not know how to calculate, namely *that the induced emission probability and the absorption probability must be equal*. This is interesting. And furthermore, in order for (42.17) and (42.12) to agree,

$$A_{mn}/B_{mn} \quad \text{must be} \quad \hbar\omega^3/\pi^2 c^2. \quad (42.18)$$

So if we know, for instance, the absorption rate for a given level, we can deduce the spontaneous emission rate and the induced emission rate, or any combination.

This is as far as Einstein or anyone else could go using such arguments. To actually compute the absolute spontaneous emission rate or the other rates for any specific atomic transition, of course, requires a knowledge of the machinery of the atom, called quantum electrodynamics, which was not discovered until eleven years later. This work of Einstein was done in 1916.

The possibility of induced emission has, today, found interesting applications. If there is light present, it will tend to induce the downward transition. The transition then adds its $\hbar\omega$ to the available light energy, if there were some atoms sitting in the upper state. Now we can arrange, by some nonthermal method, to have a gas in which the number in the state m is very much greater than the number in the state n . This is far out of equilibrium, and so is not given by the formula $e^{-\hbar\omega/kT}$, which is for equilibrium. We can even arrange it so that the number in the upper state is very large, while the number in the lower state is practically zero. Then light which has the frequency corresponding to the energy difference $E_m - E_n$ will not be strongly absorbed, because there are not many atoms in state n to absorb it. On the other hand, when that light is present, it will induce the emission from this upper state! So, if we had a lot of atoms in the upper state, there would be a sort of chain reaction, in which, the moment the atoms began to emit, more would be caused to emit, and the whole lot of them would dump down together. This is what is called a *laser*, or, in the case of the far infrared, a *maser*.

Various tricks can be used to obtain the atoms in state m . There may be higher levels to which the atoms can get if we shine in a strong beam of light of high frequency. From these high levels, they may trickle down, emitting various photons, until they all get stuck in the state m . If they tend to stay in the state m without emitting, the state is called *metastable*. And then they are all dumped down together by induced emissions. One more technical point—if we put this system in an ordinary box, it would radiate in so many different directions spontaneously, compared with the induced effect, that we would still be in trouble. But we can enhance the induced effect, increase its efficiency, by putting nearly perfect mirrors on each side of the box, so that the light which is emitted gets another chance, and another chance, and another chance, to induce more emission. Although the mirrors are almost one hundred percent reflecting, there is a slight amount of transmission of the mirror, and a little light gets out. In the end, of course, from the conservation of energy, all the light goes out in a nice uniform straight direction which makes the strong light beams that are possible today with lasers.

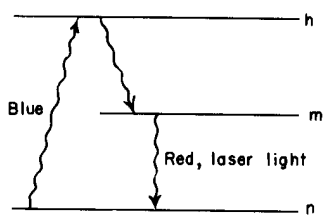


Fig. 42-3. By exciting, say by blue light, a higher state h , which may emit a photon leaving atoms in state m , the number in this state m becomes sufficiently large to start laser action.

Diffusion

43-1 Collisions between molecules

We have considered so far only the molecular motions in a gas which is in thermal equilibrium. We want now to discuss what happens when things are near, but not exactly in, equilibrium. In a situation far from equilibrium, things are extremely complicated, but in a situation very close to equilibrium we can easily work out what happens. To see what happens, we must, however, return to the kinetic theory. Statistical mechanics and thermodynamics deal with the equilibrium situation, but away from equilibrium we can only analyze what occurs atom by atom, so to speak.

As a simple example of a nonequilibrium circumstance, we shall consider the diffusion of ions in a gas. Suppose that in a gas there is a relatively small concentration of ions—electrically charged molecules. If we put an electric field on the gas, then each ion will have a force on it which is different from the forces on the neutral molecules of the gas. If there were no other molecules present, an ion would have a constant acceleration until it reached the wall of the container. But because of the presence of the other molecules, it cannot do that; its velocity increases only until it collides with a molecule and loses its momentum. It starts again to pick up more speed, but then it loses its momentum again. The net effect is that an ion works its way along an erratic path, but with a net motion in the direction of the electric force. We shall see that the ion has an average “drift” with a mean speed which is proportional to the electric field—the stronger the field, the faster it goes. While the field is on, and while the ion is moving along, it is, of course, *not* in thermal equilibrium, it is trying to get to equilibrium, which is to be sitting at the end of the container. By means of the kinetic theory we can compute the drift velocity.

It turns out that with our present mathematical abilities we cannot really compute *precisely* what will happen, but we can obtain approximate results which exhibit all the essential features. We can find out how things will vary with pressure, with temperature, and so on, but it will not be possible to get precisely the correct numerical factors in front of all the terms. We shall, therefore, in our derivations, not worry about the precise value of numerical factors. They can be obtained only by a very much more sophisticated mathematical treatment.

Before we consider what happens in nonequilibrium situations, we shall need to look a little closer at what goes on in a gas in thermal equilibrium. We shall need to know, for example, what the average time between successive collisions of a molecule is.

Any molecule experiences a sequence of collisions with other molecules—in a random way, of course. A particular molecule will, in a long period of time T , have a certain number, N , of hits. If we double the length of time, there will be twice as many hits. So the number of collisions is proportional to the time T . We would like to write it this way:

$$N = T/\tau \quad (43.1)$$

We have written the constant of proportionality as $1/\tau$, where τ will have the dimensions of a time. The constant τ is the average time between collisions. Suppose, for example, that in an hour there are 60 collisions; then τ is one minute. We would say that τ (one minute) is the *average time* between the collisions.

We may often wish to ask the following question: “What is the *chance* that a molecule will experience a collision during the next *small interval* of time dt ?” The answer, we may intuitively understand, is dt/τ . But let us try to make a more

43-1 Collisions between molecules

43-2 The mean free path

43-3 The drift speed

43-4 Ionic conductivity

43-5 Molecular diffusion

43-6 Thermal conductivity

convincing argument. Suppose that there were a very large number N of molecules. How many will have collisions in the next interval of time dt ? If there is equilibrium, nothing is changing *on the average* with time. So N molecules waiting the time dt will have the same number of collisions as *one* molecule waiting for the time $N dt$. That number we know is $N dt/\tau$. So the number of hits of N molecules is $N dt/\tau$ in a time dt , and the chance, or probability, of a hit for any one molecule is just $1/N$ as large, or $(1/N)(N dt/\tau) = dt/\tau$, as we guessed above. That is to say, the fraction of the molecules which will suffer a collision in the time dt is dt/τ . To take an example, if τ is one minute, then in one second the fraction of particles which will suffer collisions is $1/60$. What this means, of course, is that $1/60$ of the molecules happen to be close enough to what they are going to hit next that *their* collisions will occur in the next minute.

When we say that τ , the mean time between collisions, is one minute, we do not mean that all the collisions will occur at times separated by exactly one minute. A particular particle does not have a collision, wait one minute, and then have another collision. The times between successive collisions are quite variable. We will not need it for our later work here, but we may make a small diversion to answer the question: "What *are* the times between collisions?" We know that for the case above, the *average* time is one minute, but we might like to know, for example, what is the chance that we get no collision for *two* minutes?

We shall find the answer to the general question: "What is the probability that a molecule will go for a time t without having a collision?" At some arbitrary instant—that we call $t = 0$ —we begin to watch a particular molecule. What is the chance that it gets by until t without colliding with another molecule? To compute the probability, we observe what is happening to all N_0 molecules in a container. After we have waited a time t , some of them will have had collisions. We let $N(t)$ be the number that have *not* had collisions up to the time t . $N(t)$ is, of course, less than N_0 . We can find $N(t)$ because we know how it changes with time. If we know that $N(t)$ molecules have got by until t , then $N(t + dt)$, the number which get by until $t + dt$, is *less* than $N(t)$ by the number that have collisions in dt . The number that collide in dt we have written above in terms of the mean time τ as $dN = N(t) dt/\tau$. We have the equation

$$N(t + dt) = N(t) - N(t) \frac{dt}{\tau}. \quad (43.2)$$

The quantity on the left-hand side, $N(t + dt)$, can be written, according to the definitions of calculus, as $N(t) + (dN/dt) dt$. Making this substitution, Eq. (43.2) yields

$$\frac{dN(t)}{dt} = - \frac{N(t)}{\tau}. \quad (43.3)$$

The number that are being lost in the interval dt is proportional to the number that are present, and inversely proportional to the mean life τ . Equation (43.3) is easily integrated if we rewrite it as

$$\frac{dN(t)}{N(t)} = - \frac{dt}{\tau}. \quad (43.4)$$

Each side is a perfect differential, so the integral is

$$\ln N(t) = -t/\tau + (\text{a constant}), \quad (43.5)$$

which says the same thing as

$$N(t) = (\text{constant})e^{-t/\tau}. \quad (43.6)$$

We know that the constant must be just N_0 , the total number of molecules present, since all of them start at $t = 0$ to wait for their "next" collision. We can write our result as

$$N(t) = N_0 e^{-t/\tau}. \quad (43.7)$$

If we wish the *probability* of no collision, $P(t)$, we can get it by dividing $N(t)$ by N_0 , so

$$P(t) = e^{-t/\tau}. \quad (43.8)$$

Our result is: the probability that a particular molecule survives a time t without a collision is $e^{-t/\tau}$, where τ is the mean time between collisions. The probability starts out at 1 (or certainty) for $t = 0$, and gets less as t gets bigger and bigger. The probability that the molecule avoids a collision for a time equal to τ is $e^{-1} = 0.37 \dots$. The chance is less than one-half that it will have a greater than average time between collisions. That is all right, because there are enough molecules which go collision-free for times much *longer* than the mean time before colliding, so that the average time can still be τ .

We originally defined τ as the average time *between* collisions. The result we have obtained in Eq. (43.7) also says that the mean time from an *arbitrary* starting instant to the *next* collision is *also* τ . We can demonstrate this somewhat surprising fact in the following way. The number of molecules which experience their *next* collision in the interval dt at the time t after an arbitrarily chosen starting time is $N(t) dt/\tau$. Their "time until the next collision" is, of course, just t . The "average time until the next collision" is obtained in the usual way:

$$\text{Average time until the next collision} = \frac{1}{N_0} \int_0^\infty t \frac{N(t) dt}{\tau}.$$

Using $N(t)$ obtained in (43.7) and evaluating the integral, we find indeed that τ is the average time from *any* instant until the next collision.

43-2 The mean free path

Another way of describing the molecular collisions is to talk not about the *time* between collisions, but about *how far* the particle moves between collisions. If we say that the average time between collisions is τ , and that the molecules have a mean velocity v , we can expect that the average *distance* between collisions, which we shall call l , is just the product of τ and v . This distance between collisions is usually called the *mean free path*:

$$\text{Mean free path } l = \tau v. \quad (43.9)$$

In this chapter we shall be a little careless about *what kind of average* we mean in any particular case. The various possible averages—the mean, the root-mean-square, etc.—are all nearly equal and differ by factors which are near to one. Since a detailed analysis is required to obtain the correct numerical factors anyway, we need not worry about which average is required at any particular point. We may also warn the reader that the algebraic symbols we are using for some of the physical quantities (e.g., l for the mean free path) do not follow a generally accepted convention, mainly because there is no general agreement.

Just as the chance that a molecule will have a collision in a short time dt is equal to dt/τ , the chance that it will have a collision in going a distance dx is dx/l . Following the same line of argument used above, the reader can show that the probability that a molecule will go at least the distance x before having its next collision is $e^{-x/l}$.

The average distance a molecule goes before colliding with another molecule—the mean free path l —will depend on how many molecules there are around and on the "size" of the molecules, i.e., how big a target they represent. The effective "size" of a target in a collision we usually describe by a "collision cross section," the same idea that is used in nuclear physics, or in light-scattering problems.

Consider a moving particle which travels a distance dx through a gas which has n_0 scatterers (molecules) per unit volume (Fig. 43-1). If we look at each unit of area perpendicular to the direction of motion of our selected particle, we will find there $n_0 dx$ molecules. If each one presents an effective collision area or, as it is usually called, "collision cross section," σ_c , then the total area covered by the scatterers is $\sigma_c n_0 dx$.

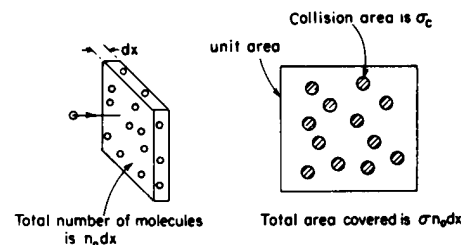


Fig. 43-1. Collision cross section.

By "collision cross section" we mean the area within which the center of our particle must be located if it is to collide with a particular molecule. If molecules were little spheres (a classical picture) we would expect that $\sigma_c = \pi(r_1 + r_2)^2$, where r_1 and r_2 are the radii of the two colliding objects. The chance that our particle will have a collision is the ratio of the area covered by scattering molecules to the total area, which we have taken to be one. So the probability of a collision in going a distance dx is just $\sigma_c n_0 dx$:

$$\text{Chance of a collision in } dx = \sigma_c n_0 dx. \quad (43.10)$$

We have seen above that the chance of a collision in dx can also be written in terms of the mean free path l as dx/l . Comparing this with (43.10), we can relate the mean free path to the collision cross section:

$$\frac{1}{l} = \sigma_c n_0, \quad (43.11)$$

which is easier to remember if we write it as

$$\sigma_c n_0 l = 1. \quad (43.12)$$

This formula can be thought of as saying that there should be one collision, on the average, when the particle goes through a distance l in which the scattering molecules *could* just cover the total area. In a cylindrical volume of length l and a base of unit area, there are $n_0 l$ scatterers; if each one has an area σ_c the total area covered is $n_0 l \sigma_c$, which is just one unit of area. The whole area is *not* covered, of course, because some molecules are partly hidden behind others. That is why some molecules go farther than l before having a collision. It is only *on the average* that the molecules have a collision by the time they go the distance l . From measurements of the mean free path l we can determine the scattering cross section σ_c , and compare the result with calculations based on a detailed theory of atomic structure. But that is a different subject! So we return to the problem of nonequilibrium states.

43-3 The drift speed

We want to describe what happens to a molecule, or several molecules, which are different in some way from the large majority of the molecules in a gas. We shall refer to the "majority" molecules as the "background" molecules, and we shall call the molecules which are different from the background molecules "special" molecules or, for short, the *S*-molecules. A molecule could be special for any number of reasons: It might be heavier than the background molecules. It might be a different chemical. It might have an electric charge—i.e., be an ion in a background of uncharged molecules. Because of their different masses or charges the *S*-molecules may have forces on them which are different from the forces on the background molecules. By considering what happens to these *S*-molecules we can understand the basic effects which come into play in a similar way in many different phenomena. To list a few: the diffusion of gases, electric currents in batteries, sedimentation, centrifugal separation, etc.

We begin by concentrating on the basic process: an *S*-molecule in a background gas is acted on by some specific force **F** (which might be, e.g., gravitational or electrical) and *in addition* by the not-so-specific forces due to collisions with the background molecules. We would like to describe the *general* behavior of the *S*-molecule. What happens to it, *in detail*, is that it darts around hither and yon as it collides over and over again with other molecules. But if we watch it carefully we see that it does make some net progress in the direction of the force **F**. We say that there is a *drift* superposed on its random motion. We would like to know what the speed of its drift is—its *drift velocity*—due to the force **F**.

If we start to observe an *S*-molecule at some instant we may expect that it is somewhere between two collisions. In addition to the velocity it was left with after its last collision it is picking up some velocity component due to the force **F**. In a

short time (on the average, in a time τ) it will experience a collision and start out on a new piece of its trajectory. It will have a new starting velocity, but the same acceleration from \mathbf{F} .

To keep things simple for the moment, we shall suppose that after each collision our S -molecule gets a completely “fresh” start. That is, that it keeps no remembrance of its past acceleration by \mathbf{F} . This might be a reasonable assumption if our S -molecule were much lighter than the background molecules, but it is certainly not valid in general. We shall discuss later an improved assumption.

For the moment, then, our assumption is that the S -molecule leaves each collision with a velocity which may be in any direction with equal likelihood. The starting velocity will take it equally in all directions and will not contribute to any net motion, so we shall not worry further about its initial velocity after a collision. In addition to its random motion, each S -molecule will have, at any moment, an additional velocity in the direction of the force \mathbf{F} , which it has picked up *since* its last collision. What is the *average* value of *this* part of the velocity? It is just the acceleration \mathbf{F}/m (where m is the mass of the S -molecule) times the *average* time *since* the last collision. Now the average time *since* the *last* collision must be the same as the average time *until* the *next* collision, which we have called τ , above. The *average* velocity from \mathbf{F} , of course, is just what is called the drift velocity, so we have the relation

$$v_{\text{drift}} = \frac{F\tau}{m}. \quad (43.13)$$

This basic relation is the heart of our subject. There may be some complication in determining what τ is, but the basic process is defined by Eq. (43.13).

You will notice that the drift velocity is *proportional* to the force. There is, unfortunately, no generally used name for the constant of proportionality. Different names have been used for each different kind of force. If in an electrical problem the force is written as the charge times the electric field, $\mathbf{F} = q\mathbf{E}$, then the constant of proportionality between the velocity and the electric field \mathbf{E} is usually called the “mobility.” In spite of the possibility of some confusion, we shall use the term *mobility* for the ratio of the drift velocity to the force for *any* force. We write

$$v_{\text{drift}} = \mu F \quad (43.14)$$

in general, and we shall call μ the mobility. We have from Eq. (43.13) that

$$\mu = \tau/m. \quad (43.15)$$

The mobility is proportional to the mean time between collisions (there are fewer collisions to slow it down) and inversely proportional to the mass (more inertia means less speed picked up between collisions).

To get the correct numerical coefficient in Eq. (43.13), which is correct as given, takes some care. Without intending to confuse, we should still point out that the arguments have a subtlety which can be appreciated only by a careful and detailed study. To illustrate that there are difficulties, in spite of appearances, we shall make over again the argument which led to Eq. (43.13) in a reasonable *but erroneous* way (and the way one will find in many textbooks!).

We might have said: The mean time between collisions is τ . After a collision the particle starts out with a random velocity, but it picks up an additional velocity between collisions, which is equal to the acceleration times the time. Since it takes the time τ to arrive at the *next* collision it gets there with the velocity $(F/m)\tau$. At the beginning of the collision it had zero velocity. So between the two collisions it has, on the average, a velocity one-half of the final velocity, so the mean drift velocity is $\frac{1}{2}F\tau/m$. (Wrong!) This result is wrong and the result in Eq. (43.13) is right, although the arguments may sound equally satisfactory. The reason the second result is wrong is somewhat subtle, and has to do with the following: The argument is made as though all collisions were separated by the mean time τ . The fact is that some times are shorter and others are longer than the mean. Short times occur *more often* but make *less* contribution to the drift velocity because they have less

chance “to really get going.” If one takes proper account of the *distribution* of free times between collisions, one can show that there should not be the factor $\frac{1}{2}$ that was obtained from the second argument. The error was made in trying to relate by a simple argument the *average final* velocity to the average velocity itself. This relationship is not simple, so it is best to concentrate on what is wanted: the average velocity itself. The first argument we gave determines the average velocity directly—and correctly! But we can perhaps see now why we shall not in general try to get all of the correct numerical coefficients in our elementary derivations!

We return now to our simplifying assumption that each collision knocks out all memory of the past motion—that a fresh start is made after each collision. Suppose our *S*-molecule is a heavy object in a background of lighter molecules. Then our *S*-molecule will not lose its “forward” momentum in each collision. It would take several collisions before its motion was “randomized” again. We should assume, instead, that at each collision—in each time τ on the average—it loses a certain fraction of its momentum. We shall not work out the details, but just state that the result is equivalent to replacing τ , the average collision time, by a new—and longer— τ which corresponds to the average “forgetting time,” i.e., the average time to forget its forward momentum. With such an interpretation of τ we can use our formula (43.15) for situations which are not quite as simple as we first assumed.

43-4 Ionic conductivity

We now apply our results to a special case. Suppose we have a gas in a vessel in which there are also some ions—atoms or molecules with a net electric charge. We show the situation schematically in Fig. 43-2. If two opposite walls of the container are metallic plates, we can connect them to the terminals of a battery and thereby produce an electric field in the gas. The electric field will result in a force on the ions, so they will begin to drift toward one or the other of the plates. An electric current will be induced, and the gas with its ions will behave like a resistor. By computing the ion flow from the drift velocity we can compute the resistance. We ask, specifically: How does the flow of electric current depend on the voltage difference V that we apply across the two plates?

We consider the case that our container is a rectangular box of length b and cross-sectional area A (Fig. 43-2). If the potential difference, or voltage, from one plate to the other is V , the electric field E between the plates is V/b . (The electric potential is the work done in carrying a unit charge from one plate to the other. The force on a unit charge is E . If E is the same everywhere between the plates, which is a good enough approximation for now, the work done on a unit charge is just Eb , so $V = Eb$.) The special force on an ion of the gas is qE , where q is the charge on the ion. The drift velocity of the ion is then μ times this force, or

$$v_{\text{drift}} = \mu F = \mu q E = \mu q \frac{V}{b}. \quad (43.16)$$

An electric current I is the flow of charge in a unit time. The electric current to one of the plates is given by the total charge of the ions which arrive at the plate in a unit of time. If the ions drift toward the plate with the velocity v_{drift} , then those which are within a distance $(v_{\text{drift}} \cdot T)$ will arrive at the plate in the time T . If there are n_i ions per unit volume, the number which reach the plate in the time T is $(n_i \cdot A \cdot v_{\text{drift}} \cdot T)$. Each ion carries the charge q , so we have that

$$\text{Charge collected in } T = q n_i A v_{\text{drift}} T. \quad (43.17)$$

The current I is the charge collected in T divided by T , so

$$I = q n_i A v_{\text{drift}}. \quad (43.18)$$

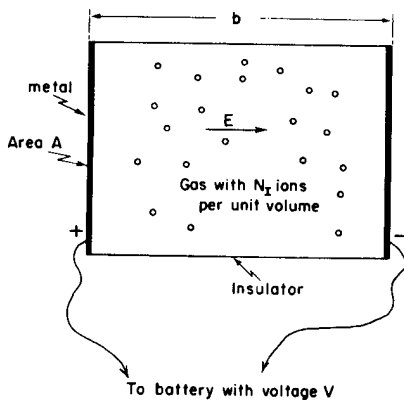


Fig. 43-2. Electric current from an ionized gas.

Substituting v_{drift} from (43.16), we have

$$I = \mu q^2 n_i \frac{A}{b} V. \quad (43.19)$$

We find that the current is proportional to the voltage, which is just the form of Ohm's law, and the resistance R is the inverse of the proportionality constant:

$$\frac{1}{R} = \mu q^2 n_i \frac{A}{b}. \quad (43.20)$$

We have a relation between the resistance and the molecular properties n_i , q , and μ , which depends in turn on m and τ . If we know n_i and q from atomic measurements, a measurement of R could be used to determine μ , and from μ also τ .

43-5 Molecular diffusion

We turn now to a different kind of problem, and a different kind of analysis: the theory of diffusion. Suppose that we have a container of gas in thermal equilibrium, and that we introduce a small amount of a different kind of gas at some place in the container. We shall call the original gas the "background" gas and the new one the "special" gas. The special gas will start to spread out through the whole container, but it will spread slowly because of the presence of the background gas. This slow spreading-out process is called *diffusion*. The diffusion is controlled mainly by the molecules of the special gas getting knocked about by the molecules of the background gas. After a large number of collisions, the special molecules end up spread out more or less evenly throughout the whole volume. We must be careful *not* to confuse diffusion of a gas with the gross transport that may occur due to convection currents. Most commonly, the mixing of two gases occurs by a combination of convection and diffusion. We are interested now only in the case that there are *no* "wind" currents. The gas is spreading only by molecular motions, by diffusion. We wish to compute how fast diffusion takes place.

We now compute the *net flow* of molecules of the "special" gas due to the molecular motions. There will be a net flow only when there is some nonuniform distribution of the molecules, otherwise all of the molecular motions would average to give no net flow. Let us consider first the flow in the x -direction. To find the flow, we consider an imaginary plane surface perpendicular to the x -axis and count the number of special molecules that cross this plane. To obtain the net flow, we must count as positive those molecules which cross in the direction of positive x and *subtract* from this number the number which cross in the negative x -direction. As we have seen many times, the number which cross a surface area in a time ΔT is given by the number which start the interval ΔT in a volume which extends the distance $v \Delta T$ from the plane. (Note that v , here, is the actual molecular velocity, not the drift velocity.)

We shall simplify our algebra by giving our surface one unit of area. Then the number of special molecules which pass from left to right (taking the $+x$ -direction to the right) is $n_- v \Delta T$, where n_- is the number of special molecules per unit volume to the left (within a factor of 2 or so, but we are ignoring such factors!). The number which cross from right to left is, similarly, $n_+ v \Delta T$, where n_+ is the number density of special molecules on the right-hand side of the plane. If we call the molecular current J , by which we mean the net flow of molecules per unit area per unit time, we have

$$J = \frac{n_- v \Delta T - n_+ v \Delta T}{\Delta T}, \quad (43.21)$$

or

$$J = (n_- - n_+)v. \quad (43.22)$$

What shall we use for n_- and n_+ ? When we say "the density on the left," how *far* to the left do we mean? We should choose the density at the place from which the molecules started their "flight," because the number which *start* such

trips is determined by the number present at that place. So by n_- we should mean the density a distance to the left equal to the mean free path l , and by n_+ , the density at the distance l to the right of our imaginary surface.

It is convenient to consider that the distribution of our special molecules in space is described by a continuous function of x , y , and z which we shall call n_a . By $n_a(x, y, z)$ we mean the number density of special molecules in a small volume element centered on (x, y, z) . In terms of n_a we can express the difference ($n_+ - n_-$) as

$$(n_+ - n_-) = \frac{dn_a}{dx} \Delta x = \frac{dn_a}{dx} \cdot 2l. \quad (43.23)$$

Substituting this result in Eq. (43.22) and neglecting the factor of 2, we get

$$J_x = -lv \frac{dn_a}{dx}. \quad (43.24)$$

We have found that the flow of special molecules is proportional to the derivative of the density, or to what is sometimes called the “gradient” of the density.

It is clear that we have made several rough approximations. Besides various factors of two we have left out, we have used v where we should have used v_x , and we have assumed that n_+ and n_- refer to places at the perpendicular distance l from our surface, whereas for those molecules which do not travel perpendicular to the surface element, l should correspond to the *slant* distance from the surface. All of these refinements can be made; the result of a more careful analysis shows that the right-hand side of Eq. (43.24) should be multiplied by $1/3$. So a better answer is

$$J_x = -\frac{lv}{3} \frac{dn_a}{dx}. \quad (43.25)$$

Similar equations can be written for the currents in the y - and z -directions.

The current J_x and the density gradient dn_a/dx can be measured by macroscopic observations. Their experimentally determined ratio is called the “diffusion coefficient,” D . That is,

$$J_x = -D \frac{dn_a}{dx}. \quad (43.26)$$

We have been able to show that for a gas we expect

$$D = \frac{1}{3}lv. \quad (43.27)$$

So far in this chapter we have considered two distinct processes: *mobility*, the drift of molecules due to “outside” forces; and *diffusion*, the spreading determined only by the internal forces, the random collisions. There is, however, a relation between them, since they both depend basically on the thermal motions, and the mean free path l appears in both calculations.

If, in Eq. (43.25), we substitute $l = v\tau$ and $\tau = \mu m$, we have

$$J_x = -\frac{1}{3}mv^2\mu \frac{dn_a}{dx}. \quad (43.28)$$

But mv^2 depends only on the temperature. We recall that

$$\frac{1}{2}mv^2 = \frac{3}{2}kT, \quad (43.29)$$

so

$$J_x = -\mu kT \frac{dn_a}{dx}. \quad (43.30)$$

We find that D , the *diffusion* coefficient, is just kT times μ , the *mobility* coefficient:

$$D = \mu kT. \quad (43.31)$$

And it turns out that the numerical coefficient in (43.31) is exactly right—no extra factors have to be thrown in to adjust for our rough assumptions. We can show,

in fact, that (43.31) must *always* be correct—even in complicated situations (for example, the case of a suspension in a liquid) where the details of our simple calculations would not apply at all.

To show that (43.31) must be correct in general, we shall derive it in a different way, using only our basic principles of statistical mechanics. Imagine a situation in which there is a gradient of “special” molecules, and we have a diffusion current proportional to the density gradient, according to Eq. (43.26). We now apply a force field in the x -direction, so that each special molecule feels the force F . According to the *definition* of the mobility μ there will be a drift velocity given by

$$v_{\text{drift}} = \mu F. \quad (43.32)$$

By our usual arguments, the *drift current* (the *net* number of molecules which pass a unit of area in a unit of time) will be

$$J_{\text{drift}} = n_a v_{\text{drift}}, \quad (43.33)$$

or

$$J_{\text{drift}} = n_a \mu F. \quad (43.34)$$

We now *adjust* the force F so that the drift current due to F just *balances* the diffusion, so that there is *no net flow* of our special molecules. We have $J_x + J_{\text{drift}} = 0$, or

$$D \frac{dn_a}{dx} = n_a \mu F. \quad (43.35)$$

Under the “balance” conditions we find a steady (with time) gradient of density given by

$$\frac{dn_a}{dx} = \frac{n_a \mu F}{D}. \quad (43.36)$$

But notice! We are describing an *equilibrium* condition, so our *equilibrium* laws of statistical mechanics apply. According to these laws the probability of finding a molecule at the coordinate x is proportional to $e^{-U/kT}$, where U is the potential energy. In terms of the number density n_a , this means that

$$n_a = n_0 e^{-U/kT}. \quad (43.37)$$

If we differentiate (43.37) with respect to x , we find

$$\frac{dn_a}{dx} = -n_0 e^{-U/kT} \cdot \frac{1}{kT} \frac{dU}{dx}, \quad (43.38)$$

or

$$\frac{dn_a}{dx} = -\frac{n_a}{kT} \frac{dU}{dx}. \quad (43.39)$$

In our situation, since the force F is in the x -direction, the potential energy U is just $-Fx$, and $-dU/dx = F$. Equation (43.39) then gives

$$\frac{dn_a}{dx} = \frac{n_a F}{kT}. \quad (43.40)$$

[This is just exactly Eq. (40.2), from which we deduced $e^{-U/kT}$ in the first place, so we have come in a circle]. Comparing (43.40) with (43.36), we get exactly Eq. (43.31). We have shown that Eq. (43.31), which gives the diffusion current in terms of the mobility, has the correct coefficient and is very generally true. Mobility and diffusion are intimately connected. This relation was first deduced by Einstein.

43-6 Thermal conductivity

The methods of the kinetic theory that we have been using above can be used also to compute the *thermal conductivity* of a gas. If the gas at the top of a container is hotter than the gas at the bottom, heat will flow from the top to the bottom. (We think of the top being hotter because otherwise convection currents would be

set up and the problem would no longer be one of heat *conduction*.) The transfer of heat from the hotter gas to the colder gas is by the diffusion of the “hot” molecules—those with more energy—downward and the diffusion of the “cold” molecules upward. To compute the flow of thermal energy we can ask about the energy carried downward across an element of area by the downward-moving molecules, and about the energy carried upward across the surface by the upward-moving molecules. The difference will give us the net downward flow of energy.

The thermal conductivity κ is defined as the ratio of the rate at which thermal energy is carried across a unit surface area, to the temperature gradient:

$$\frac{1}{A} \frac{dQ}{dt} = -\kappa \frac{dT}{dz}. \quad (43.41)$$

Since the details of the calculations are quite similar to those we have done above in considering the flow of electric current in an ionized gas, we shall leave it as an exercise for the reader to show that

$$\kappa = \frac{kn\bar{v}}{\gamma - 1}, \quad (43.42)$$

where $(\gamma - 1)kT$ is the average energy of a molecule at the temperature T .

If we use our relation $n\bar{\sigma}_c = 1$, the heat conductivity can be written as

$$\kappa = \frac{1}{\gamma - 1} \frac{kv}{\bar{\sigma}_c}. \quad (43.43)$$

We have a rather surprising result. We know that the average velocity of gas molecules depends on the temperature but *not on the density*. We expect $\bar{\sigma}_c$ to depend only on the *size* of the molecules. So our simple result says that the thermal conductivity κ (and therefore the *rate* of flow of heat in any particular circumstance) is independent of the *density* of the gas! The change in the number of “carriers” of energy with a change in density is just compensated by the larger distance the “carriers” can go between collisions.

One may ask: “Is the heat flow independent of the gas density in the limit as the density goes to zero? When there is no gas at all?” Certainly not! The formula (43.43) was derived, as were all the others in this chapter, under the assumption that the mean free path between collisions is much smaller than any of the dimensions of the container. Whenever the gas density is so low that a molecule has a fair chance of crossing from one wall of its container to the other without having a collision, none of the calculations of this chapter apply. We must in such cases go back to kinetic theory and calculate again the details of what will occur.

The Laws of Thermodynamics

44-1 Heat engines; the first law

So far we have been discussing the properties of matter from the atomic point of view, trying to understand roughly what will happen if we suppose that things are made of atoms obeying certain laws. However, there are a number of relationships among the properties of substances which can be worked out without consideration of the detailed structure of the materials. The determination of the relationships among the various properties of materials, without knowing their internal structure, is the subject of *thermodynamics*. Historically, thermodynamics was developed before an understanding of the internal structure of matter was achieved.

To give an example: we know from the kinetic theory that the pressure of a gas is caused by molecular bombardment, and we know that if we heat a gas, so that the bombardment increases, the pressure must increase. Conversely, if the piston in a container of the gas is moved inward against the force of bombardment, the energy of the molecules bombarding the piston will increase, and consequently the temperature will increase. So, on the one hand, if we increase the temperature at a given volume, we increase the pressure. On the other hand, if we compress the gas, we will find that the temperature will rise. From the kinetic theory, one can derive a quantitative relationship between these two effects, but instinctively one might guess that they are related in some necessary fashion which is independent of the details of the collisions.

Let us consider another example. Many people are familiar with this interesting property of rubber: If we take a rubber band and pull it, it gets warm. If one puts it between his lips, for example, and pulls it out, he can feel a distinct warming, and this warming is reversible in the sense that if he relaxes the rubber band quickly while it is between his lips, it is distinctly cooled. That means that when we stretch a rubber band it heats, and when we release the tension of the band it cools. Now our instincts might suggest that if we heated a band, it might pull: that the fact that pulling a band heats it might imply that heating a band should cause it to contract. And, in fact, if we apply a gas flame to a rubber band holding a weight, we will see that the band contracts abruptly (Fig. 44-1). So it is true that when we heat a rubber band it pulls, and this fact is definitely related to the fact that when we release the tension of it, it cools.

The internal machinery of rubber that causes these effects is quite complicated. We will describe it from a molecular point of view to some extent, although our main purpose in this chapter is to understand the relationship of these effects independently of the molecular model. Nevertheless, we can show from the molecular model that the effects are closely related. One way to understand the behavior of rubber is to recognize that this substance consists of an enormous tangle of long chains of molecules, a kind of "molecular spaghetti," with one extra complication: between the chains there are cross-links—like spaghetti that is sometimes welded together where it crosses another piece of spaghetti—a grand tangle. When we pull out such a tangle, some of the chains tend to line up along the direction of the pull. At the same time, the chains are in thermal motion, so they hit each other continually. It follows that such a chain, if stretched, would not by itself remain stretched, because it would be hit from the sides by the other chains and other molecules, and would tend to kink up again. So the real reason why a rubber band tends to contract is this: when one pulls it out, the chains are lengthwise, and the thermal agitations of the molecules on the sides of the chains tend

44-1 Heat engines; the first law

44-2 The second law

44-3 Reversible engines

44-4 The efficiency of an ideal engine

44-5 The thermodynamic temperature

44-6 Entropy

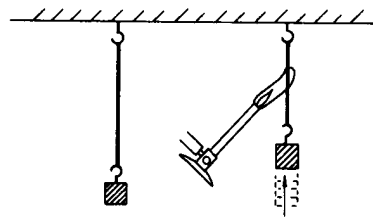


Fig. 44-1. The heated rubber band.

to kink the chains up, and make them shorten. One can then appreciate that if the chains are held stretched and the temperature is increased, so that the vigor of the bombardment on the sides of the chains is also increased, the chains tend to pull in, and they are able to pull a stronger weight when heated. If, after being stretched for a time, a rubber band is allowed to relax, each chain becomes soft, and the molecules striking it lose energy as they pound into the relaxing chain. So the temperature falls.

We have seen how these two process, contraction when heated and cooling during relaxation, can be related by the kinetic theory, but it would be a tremendous challenge to determine from the theory the precise relationship between the two. We would have to know how many collisions there were each second and what the chains look like, and we would have to take account of all kinds of other complications. The detailed mechanism is so complex that we cannot, by kinetic theory, really determine exactly what happens; still, a definite relation between the two effects we observe can be worked out without knowing anything about the internal machinery!

The whole subject of thermodynamics depends essentially upon the following kind of consideration: because a rubber band is "stronger" at higher temperatures than it is at lower temperatures, it ought to be possible to lift weights, and to move them around, and thus to do work with heat. In fact, we have already seen experimentally that a heated rubber band can lift a weight. The study of the way that one does work with heat is the beginning of the science of thermodynamics. Can we make an engine which uses the heating effect on a rubber band to do work? One can make a silly looking engine that does just this. It consists of a bicycle wheel in which all the spokes are rubber bands (Fig. 44-2). If one heats the rubber bands on one side of the wheel with a pair of heat lamps, they become "stronger" than the rubber bands on the other side. The center of gravity of the wheel will be pulled to one side, away from the bearing, so that the wheel turns. As it turns, cool rubber bands move toward the heat, and the heated bands move away from the heat and cool, so that the wheel turns slowly so long as the heat is applied. The efficiency of this engine is extremely low. Four hundred watts of power pour into the two lamps, but it is just possible to lift a fly with such an engine! An interesting question, however, is whether we can get heat to do the work in more efficient ways.

In fact, the science of thermodynamics began with an analysis, by the great engineer Sadi Carnot, of the problem of how to build the best and most efficient engine, and this constitutes one of the few famous cases in which engineering has contributed fundamentally to physical theory. Another example that comes to mind is the more recent analysis of information theory by Claude Shannon. These two analyses, incidentally, turn out to be closely related.

Now the way a steam engine ordinarily operates is that heat from a fire boils some water, and the steam so formed expands and pushes on a piston which makes a wheel go around. So the steam pushes the piston—what then? One has to finish the job: a stupid way to complete the cycle would be to let the steam escape into the air, for then one has to keep supplying water. It is cheaper—more efficient—to let the steam go into another box, where it is condensed by cool water, and then pump the water back into the boiler, so that it circulates continuously. Heat is thus supplied to the engine and converted into work. Now would it be better to use alcohol? What property should a substance have so that it makes the best possible engine? That was the question to which Carnot addressed himself, and one of the by-products was the discovery of the type of relationship that we have just explained above.

The results of thermodynamics are all contained implicitly in certain apparently simple statements called the *laws of thermodynamics*. At the time when Carnot lived, the first law of thermodynamics, the conservation of energy, was not known. Carnot's arguments were so carefully drawn, however, that they are valid even though the first law was not known in his time! Some time afterwards, Clausius made a simpler derivation that could be understood more easily than Carnot's very subtle reasoning. But it turned out that Clausius assumed, not the conserva-

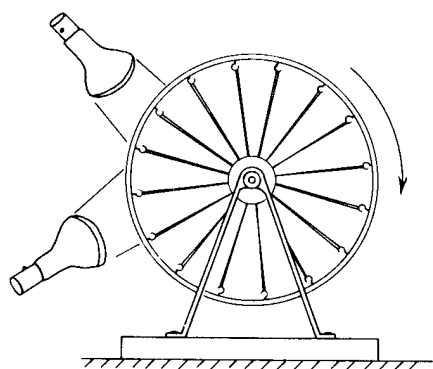


Fig. 44-2. The rubber-band heat engine.

tion of energy in general, but that *heat* was conserved according to the caloric theory, which was later shown to be false. So it has often been said that Carnot's logic was wrong. But his logic was quite correct. Only Clausius' simplified version, that everybody read, was incorrect.

The so-called second law of thermodynamics was thus discovered by Carnot before the first law! It would be interesting to give Carnot's argument that did not use the first law, but we shall not do so because we want to learn physics, not history. We shall use the first law from the start, in spite of the fact that a great deal can be done without it.

Let us begin by stating the first law, the conservation of energy: if one has a system and puts heat into it, and does work on it, then its energy is increased by the heat put in and the work done. We can write this as follows: The heat Q put into the system, plus the work W done on the system, is the increase in the energy U of the system; the latter energy is sometimes called the internal energy:

$$\text{Change in } U = Q + W. \quad (44.1)$$

The change in U can be represented as adding a little heat ΔQ and adding a little work ΔW :

$$\Delta U = \Delta Q + \Delta W, \quad (44.2)$$

which is a differential form of the same law. We know that very well, from an earlier chapter.

44-2 The second law

Now, what about the second law of thermodynamics? We know that if we do work against friction, say, the work lost to us is equal to the heat produced. If we do work in a room at temperature T , and we do the work slowly enough, the room temperature does not change much, and we have converted work into heat at a given temperature. What about the reverse possibility? Is it possible to convert the heat back into work at a given temperature? The second law of thermodynamics asserts that it is not. It would be very convenient to be able to convert heat into work merely by reversing a process like friction. If we consider only the conservation of energy, we might think that heat energy, such as that in the vibrational motions of molecules, might provide a goodly supply of useful energy. But Carnot assumed that it is impossible to extract the energy of heat at a single temperature. In other words, if the whole world were at the same temperature, one could not convert any of its heat energy into work: while the process of making work go into heat can take place at a given temperature, one cannot reverse it to get the work back again. Specifically, Carnot assumed that heat cannot be taken in at a certain temperature and converted into work *with no other change* in the system or the surroundings.

That last phrase is very important. Suppose we have a can of compressed air at a certain temperature, and we let the air expand. It can do work; it can make hammers go, for example. It cools off a little in the expansion, but if we had a big sea, like the ocean, at a given temperature—a heat reservoir—we could warm it up again. So we have taken the heat out of the sea, and we have done work with the compressed air. But Carnot was not wrong, because *we did not leave everything as it was*. If we recompress the air that we let expand, we will find we are doing extra work, and when we are finished we will discover that we not only got no work out of the system at temperature T , but we actually put some in. We must talk only about situations in which the *net result* of the whole process is to take heat away and convert it into work, just as the net result of the process of doing work against friction is to take work and convert it into heat. If we move in a circle, we can bring the system back precisely to its starting point, with the net result that we did work against friction and produced heat. Can we reverse the process? Turn a switch, so that everything goes backwards, so the friction does work against us, and cools the sea? According to Carnot: no! So let us suppose that this is impossible.

If it were possible it would mean, among other things, that we could take heat out of a cold body and put it into a hot body at no cost, as it were. Now we know it is natural that a hot thing can warm up a cool thing; if we simply put a hot body and a cold one together, and change nothing else, our experience assures us that it is not going to happen that the hot one gets hotter, and the cold one gets colder! But if we could obtain work by extracting the heat out of the ocean, say, or from anything else at a single temperature, then that work could be converted back into heat by friction at some other temperature. For instance, the other arm of a working machine could be rubbing something that is already hot. The net result would be to take heat from a “cold” body, the ocean, and to put it into a hot body. Now, the hypothesis of Carnot, the second law of thermodynamics, is sometimes stated as follows: heat cannot, of itself, flow from a cold to a hot object. But, as we have just seen, these two statements are equivalent: first, that one cannot devise a process whose only result is to convert heat to work at a single temperature, and second, that one cannot make heat flow by itself from a cold to a hot place. We shall mostly use the first form.

Carnot’s analysis of heat engines is quite similar to the argument that we gave about weight-lifting engines in our discussion of the conservation of energy in Chapter 4. In fact, that argument was patterned after Carnot’s argument about heat engines, and so the present treatment will sound very much the same.

Suppose we build a heat engine that has a “boiler” somewhere at a temperature T_1 . A certain heat Q_1 is taken from the boiler, the steam engine does some work W , and it then delivers some heat Q_2 into a “condenser” at another temperature T_2 (Fig. 44-3). Carnot did not say how much heat, because he did not know the first law, and he did not use the law that Q_2 was equal to Q_1 because he did not believe it. Although everybody thought that, according to the caloric theory, the heats Q_1 and Q_2 would have to be the same, Carnot did not say they were the same—that is part of the cleverness of his argument. If we do use the first law, we find that the heat delivered, Q_2 , is the heat Q_1 that was put in minus the work W that was done:

$$Q_2 = Q_1 - W. \quad (44.3)$$

(If we have some kind of cyclic process where water is pumped back into the boiler after it is condensed, we will say that we have heat Q_1 absorbed and work W done, during each cycle, for a certain amount of water that goes around the cycle.)

Now we shall build another engine, and see if we cannot get more work from the same amount of heat being delivered at the temperature T_1 , with the condenser still at the temperature T_2 . We shall use the same amount of heat Q_1 from the boiler, and we shall try to get more work than we did out of the steam engine, perhaps by using another fluid, such as alcohol.

44-3 Reversible engines

Now we must analyze our engines. One thing is clear: we will lose something if the engines contain devices in which there is friction. The best engine will be a frictionless engine. We assume, then, the same idealization that we did when we studied the conservation of energy; that is, a perfectly frictionless engine.

We must also consider the analog of frictionless motion, “frictionless” heat transfer. If we put a hot object at a high temperature against a cold object, so that the heat flows, then it is not possible to make that heat flow in a reverse direction by a very small change in the temperature of either object. But when we have a practically frictionless machine, if we push it with a little force one way, it goes that way, and if we push it with a little force the other way, it goes the other way. We need to find the analog of frictionless motion: heat transfer whose direction we can reverse with only a tiny change. If the difference in temperature is finite, that is impossible, but if one makes sure that heat flows always between two things at essentially the same temperature, with just an infinitesimal difference to make it flow in the desired direction, the flow is said to be reversible (Fig. 44-4). If we heat

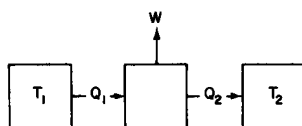


Fig. 44-3. Heat engine.

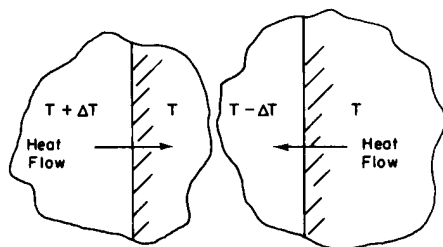


Fig. 44-4. Reversible heat transfer.

the object on the left a little, heat will flow to the right; if we cool it a little, heat will flow to the left. So we find that the ideal engine is a so-called *reversible* engine, in which every process is reversible in the sense that, by minor changes, infinitesimal changes, we can make the engine go in the opposite direction. That means that nowhere in the machine must there be any appreciable friction, and nowhere in the machine must there be any place where the heat of the reservoirs, or the flame of the boiler, is in direct contact with something definitely cooler or warmer.

Let us now consider an idealized engine in which all the processes are reversible. To show that such a thing is possible in principle, we will give an example of an engine cycle which may or may not be practical, but which is at least reversible, in the sense of Carnot's idea. Suppose that we have a gas in a cylinder equipped with a frictionless piston. The gas is not necessarily a perfect gas. The fluid does not even have to be a gas, but to be specific let us say we do have a perfect gas. Also, suppose that we have two heat pads, T_1 and T_2 —great big things that have definite temperatures, T_1 and T_2 . We will suppose in this case that T_1 is higher than T_2 . Let us first heat the gas and at the same time expand it, while it is in contact with the heat pad at T_1 . As we do this, pulling the piston out very slowly as the heat flows into the gas, we will make sure that the temperature of the gas never gets very far from T_1 . If we pull the piston out too fast, the temperature of the gas will fall too much below T_1 and then the process will not be quite reversible, but if we pull it out slowly enough, the temperature of the gas will never depart much from T_1 . On the other hand, if we push the piston back slowly, the temperature would be only infinitesimally higher than T_1 , and the heat would pour back. We see that such an isothermal (constant-temperature) expansion, done slowly and gently enough, is a reversible process.

To understand what we are doing, we shall use a plot (Fig. 44-6) of the pressure of the gas against its volume. As the gas expands, the pressure falls. The curve marked (1) tells us how the pressure and volume change if the temperature is kept fixed at the value T_1 . For an ideal gas this curve would be $PV = NkT_1$. During an isothermal expansion the pressure falls as the volume increases until we stop at the point b . At the same time, a certain heat Q_1 must flow into the gas from the reservoir, for if the gas were expanded without being in contact with the reservoir it would cool off, as we already know. Having completed the isothermal expansion, stopping at the point b , let us take the cylinder away from the reservoir and continue the expansion. This time we permit no heat to enter the cylinder. Again we perform the expansion slowly, so there is no reason why we cannot reverse it, and we again assume there is no friction. The gas continues to expand and the temperature falls, since there is no longer any heat entering the cylinder.

We let the gas expand, following the curve marked (2), until the temperature falls to T_2 , at the point marked c . This kind of expansion, made without adding heat, is called an *adiabatic* expansion. For an ideal gas, we already know that curve (2) has the form $PV^\gamma = \text{constant}$, where γ is a constant greater than 1, so that the adiabatic curve has a more negative slope than the isothermal curve. The gas cylinder has now reached the temperature T_2 , so that if we put it on the heat pad at temperature T_2 there will be no irreversible changes. Now we slowly compress the gas while it is in contact with the reservoir at T_2 , following the curve marked (3) (Fig. 44-5, Step 2). Because the cylinder is in contact with the reservoir, the temperature does not rise, but heat Q_2 flows from the cylinder into the reservoir at the temperature T_2 . Having compressed the gas isothermally along curve (3) to the point d , we remove the cylinder from the heat pad at temperature T_2 and compress it still further, without letting any heat flow out. The temperature will rise, and the pressure will follow the curve marked (4). If we carry out each step properly, we can return to the point a at temperature T_1 where we started, and repeat the cycle.

We see that on this diagram we have carried the gas around a complete cycle, and during one cycle we have put Q_1 in at temperature T_1 , and have removed Q_2 at temperature T_2 . Now the point is that this cycle is reversible, so that we could represent all the steps the other way around. We could have gone backwards instead of forwards: we could have started at point a , at temperature T_1 , expanded

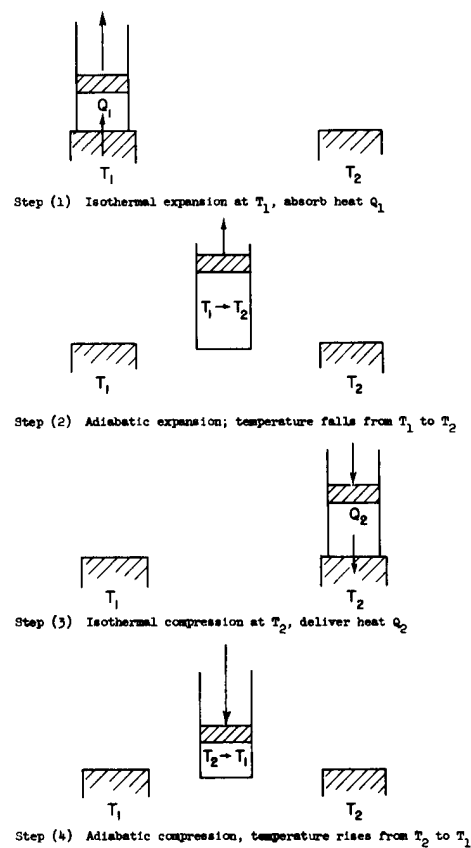


Fig. 44-5. Steps in Carnot cycle.

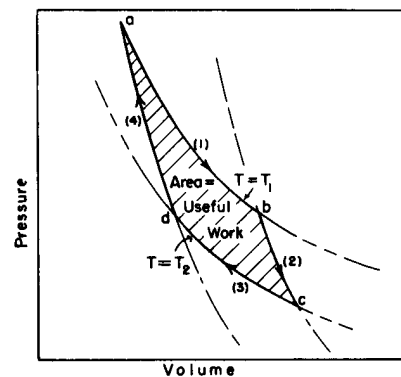


Fig. 44-6. The Carnot cycle.

along the curve (4), expanded further at the temperature T_2^* , absorbing heat Q_2 , and so on, going around the cycle backward. If we go around the cycle in one direction, we must do work on the gas; if we go in the other direction, the gas does work on us.

Incidentally, it is easy to find out what the total amount of work is, because the work during any expansion is the pressure times the change in volume, $\int P dV$. On this particular diagram, we have plotted P vertically and V horizontally. So if we call the vertical distance y and the horizontal distance x , this is $\int y dx$ —in other words, the area under the curve. So the area under each of the numbered curves is a measure of the work done by or on the gas in the corresponding step. It is easy to see that the net work done is the shaded area of the picture.

Now that we have given a single example of a reversible machine, we shall suppose that other such engines are also possible. Let us assume that we have a reversible engine A which takes Q_1 at T_1 , does work W , and delivers some heat at T_2 . Now let us assume we have any other engine B , made by man, already designed or not yet invented, made of rubber bands, steam, or whatever, reversible or not, which is designed so that it takes in the same amount of heat Q_1 at T_1 , and rejects the heat at the lower temperature T_2 (Fig. 44-7). Assume that engine B does some work, W' . Now we shall show that W' is not greater than W —that no engine can do more work than a reversible one. Why? Suppose that, indeed, W' were bigger than W . Then we could take the heat Q_1 out of the reservoir at T_1 , and with engine B we could do work W' and deliver some heat to the reservoir at T_2 ; we do not care how much. That done, we could save some of the work W' , which is supposed to be greater than W ; we could use a part of it, W , and save the remainder, $W' - W$, for useful work. With the work W we could run engine A backwards *because it is a reversible engine*. It will absorb some heat from the reservoir at T_2 and deliver Q_1 back to the reservoir at T_1 . After this double cycle, the net result would be that we would have put everything back the way it was before, and we would have done some excess work, namely $W' - W$, and *all* we would have done would be to extract energy from the reservoir at T_2 ! We were careful to restore the heat Q_1 to the reservoir at T_1 . So that reservoir can be small and “inside” our combined machine $A + B$, whose net effect is therefore to extract a net heat $W' - W$ from the reservoir at T_2 and convert it into work. But to obtain useful work from a reservoir at a single temperature *with no other changes* is impossible according to Carnot’s postulate; it cannot be done. Therefore no engine which absorbs a given amount of heat from a higher temperature T_1 and delivers it at the temperature T_2 can do more work than a reversible engine operating under the same temperature conditions.

Now suppose that engine B is also reversible. Then, of course, not only must W' be not greater than W , but now we can reverse the argument and show that W cannot be greater than W' . So, if both engines are reversible they must both do the same amount of work, and we thus come to Carnot’s brilliant conclusion: that if an engine is reversible, it makes no difference how it is designed, because the amount of work one will obtain if the engine absorbs a given amount of heat at temperature T_1 and delivers heat at some other temperature T_2 *does not depend on the design of the engine*. It is a property of the world, not a property of a particular engine.

If we could find out what the law is that determines how much work we obtain when we absorb the heat Q_1 at T_1 and deliver heat at T_2 , this quantity would be a universal thing, independent of the substance. Of course if we knew the properties of a particular substance, we could work it out and then say that all other substances must give the same amount of work in a reversible engine. That is the key idea, the clue by which we can find the relationship between how much, for instance, a rubber band contracts when we heat it, and how much it cools when we let it contract. Imagine that we put that rubber band in a reversible machine, and that we make it go around a reversible cycle. The net result, the total amount of work done, is that universal function, that great function which is independent of substance. So we see that a substance’s properties must be limited in a certain way; one

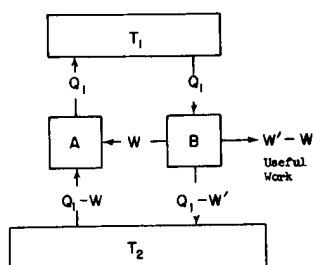


Fig. 44-7. Reversible engine A being driven backwards by engine B .

cannot make up anything he wants, or he would be able to invent a substance which he could use to produce more than the maximum allowable work when he carried it around a reversible cycle. This principle, this limitation, is the only real rule that comes out of the thermodynamics.

44-4 The efficiency of an ideal engine

Now we shall try to find the law which determines the work W as a function of Q_1 , T_1 , and T_2 . It is clear that W is proportional to Q_1 , for if we consider two reversible engines in parallel, both working together and both double engines, the combination is also a reversible engine. If each one absorbed heat Q_1 , the two together absorb $2Q_1$ and the work done is $2W$, and so on. So it is not unreasonable that W is proportional to Q_1 .

Now the next important step is to find this universal law. We can, and will, do so by studying a reversible engine with the one particular substance whose laws we know, a perfect gas. It is also possible to obtain the rule by a purely logical argument, using no particular substance at all. This is one of the very beautiful pieces of reasoning in physics and we are reluctant not to show it to you, so for those who would like to see it we shall discuss it in just a moment. But first we shall use the much less abstract and simpler method of direct calculation for a perfect gas.

We need only obtain formulas for Q_1 and Q_2 (for W is just $Q_1 - Q_2$), the heats exchanged with the reservoirs during the isothermal expansion or contraction. For example, how much heat Q_1 is absorbed from the reservoir at temperature T_1 during the isothermal expansion [marked (1) in Fig. 44-6] from point a , at pressure p_a , volume V_a , temperature T_1 , to point b with pressure p_b , volume V_b , and the same temperature T_1 ? For a perfect gas each molecule has an energy that depends only on the temperature, and since the temperature and the number of molecules are the same at a and at b , the internal energy is the same. *There is no change in U* ; all the work done by the gas,

$$W = \int_a^b p dV,$$

during the expansion is energy Q_1 taken from the reservoir. During the expansion, $pV = NkT_1$, or

$$p = \frac{NkT_1}{V}$$

or

$$Q_1 = \int_a^b p dV = \int_a^b NkT_1 \frac{dV}{V} \quad (44.4)$$

or

$$Q_1 = NkT_1 \ln \frac{V_b}{V_a}$$

is the heat taken from the reservoir at T_1 . In the same way, for the compression at T_2 [curve (3) of Fig. 44-6] the heat delivered to the reservoir at T_2 is

$$Q_2 = NkT_2 \ln \frac{V_c}{V_d}. \quad (44.5)$$

To finish our analysis we need only find a relation between V_c/V_d and V_b/V_a . This we do by noting that (2) is an adiabatic expansion from b to c , during which pV^γ is a constant. Since $pV = NkT$, we can write this as $(pV)V^{\gamma-1} = \text{const}$ or, in terms of T and V , as $TV^{\gamma-1} = \text{const}$, or

$$T_1 V_b^{\gamma-1} = T_2 V_c^{\gamma-1}. \quad (44.6)$$

Likewise, since (4), the expansion from d to a , is also adiabatic, we find

$$T_1 V_a^{\gamma-1} = T_2 V_d^{\gamma-1}. \quad (44.6a)$$

If we divide this equation by the previous one, we find that V_b/V_a must equal V_c/V_d , so the \ln 's in (44.4) and (44.5) are equal, and that

$$\frac{Q_1}{T_1} = \frac{Q_2}{T_2}. \quad (44.7)$$

This is the relation we were seeking. Although proved for a perfect gas engine, we know it must be true *for any reversible engine at all*.

Now we shall see how this universal law could also be obtained by logical argument, without knowing the properties of any specific substances, as follows. Suppose that we have three engines and three temperatures, let us say T_1 , T_2 , and T_3 . Let one engine absorb heat Q_1 from the temperature T_1 and do a certain amount of work W_{13} , and let it deliver heat Q_3 to the temperature T_3 (Fig. 44-8). Let another engine run backwards between T_2 and T_3 . Suppose that we let the second engine be of such a size that it will absorb the same heat Q_3 , and deliver the heat Q_2 . We will have to put a certain amount of work, W_{32} , into it—negative because the engine is running backwards. When the first machine goes through a cycle, it absorbs heat Q_1 and delivers Q_3 at the temperature T_3 ; then the second machine takes the same heat Q_3 out of the reservoir at the temperature T_3 and delivers it into the reservoir at temperature T_2 . Therefore the net result of the two machines in tandem is to take the heat Q_1 from T_1 , and deliver Q_2 at T_2 . The two machines are thus equivalent to a third one, which absorbs Q_1 at T_1 , does work W_{12} , and delivers heat Q_2 at T_2 , because $W_{12} = W_{13} - W_{32}$, as one can immediately show from the first law, as follows:

$$W_{13} - W_{32} = (Q_1 - Q_3) - (Q_2 - Q_3) = Q_1 - Q_2 = W_{12}. \quad (44.8)$$

We can now obtain the laws which relate the efficiencies of the engines, because there clearly must be some kind of relationship between the efficiencies of engines running between the temperatures T_1 and T_3 , and between T_2 and T_3 , and between T_1 and T_2 .

We can make the argument very clear in the following way: We have just seen that we can always relate the heat absorbed at T_1 to the heat delivered at T_2 by finding the heat delivered at some other temperature T_3 . Therefore we can get all the engines' properties if we introduce a standard temperature, analyzing everything with that standard temperature. In other words, if we knew the efficiency of an engine running between a certain temperature T and a certain arbitrary standard temperature, then we could work out the efficiency for any other difference in temperature. Because we assume we are using only reversible engines, we can work from the initial temperature down to the standard temperature and back up to the final temperature again. We shall define the standard temperature arbitrarily as *one degree*. We shall also adopt a special symbol for the heat which is delivered at this standard temperature: we shall call it Q_s . In other words, when a reversible engine absorbs the heat Q at temperature T , it will deliver, at the unit temperature, a heat Q_s . *If one engine, absorbing heat Q_1 at T_1 , delivers the heat Q_s at one degree, and if an engine absorbing heat Q_2 at temperature T_2 will also deliver the same heat Q_s at one degree, then it follows that an engine which absorbs heat Q_1 at the temperature T_1 will deliver heat Q_2 if it runs between T_1 and T_2 , as we have already proved by considering engines running between three temperatures.* So all we really have to do is to find how much heat Q_1 we need to put in at the temperature T_1 in order to deliver a certain amount of heat Q_s at the unit temperature. If we discover that, we have everything. The heat Q , of course, is a function of the temperature T . It is easy to see that the heat must increase as the temperature increases, for we know that it takes work to run an engine backwards and deliver heat at a higher temperature. It is also easy to see that the heat Q_1 must be proportional to Q_s . So the great law is something like this: for a given amount of heat Q_s delivered at one degree from an engine running at temperature T degrees, the heat Q absorbed must be that amount Q_s times some increasing function of the temperature:

$$Q = Q_s f(T). \quad (44.9)$$

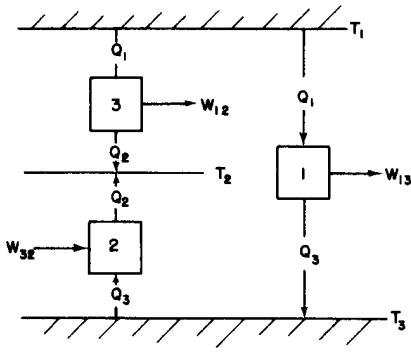


Fig. 44-8. Engines 1 and 2 together are equivalent to engine 3.

44-5 The thermodynamic temperature

At this stage we are not going to try to find the formula for the above increasing function of the temperature in terms of our familiar mercury temperature scale, but instead *we shall define temperature by a new scale*. At one time "the temperature" was defined arbitrarily by dividing the expansion of water into even degrees of a certain size. But when one then measures temperature with a mercury thermometer, one finds that the degrees are no longer even. But *now we can make a definition of temperature which is independent of any particular substance*. We can use that function $f(T)$, which does not depend on what device we use, because the efficiency of these reversible engines is independent of their working substances. Since the function we found is rising with temperature, we will define the *function itself* as the temperature, measured in units of the standard one-degree temperature, as follows:

$$Q = ST, \quad (44.10)$$

where

$$Q_S = S \cdot 1^\circ. \quad (44.11)$$

This means that we can tell how hot an object is by finding out how much heat is absorbed by a reversible engine working between the temperature of the object and the unit temperature (Fig. 44-9). If seven times more heat is taken out of a boiler than is delivered at a one-degree condenser, the temperature of the boiler will be called seven degrees, and so forth. So, by measuring how much heat is absorbed at different temperatures, we determine the temperature. The temperature defined in this way is called the *absolute thermodynamic temperature*, and it is independent of the substance. We shall use this definition exclusively from now on.*

Now we see that when we have two engines, one working between T_1 and one degree, the other working between T_2 and one degree, delivering the same heat at unit temperature, then the heats absorbed must be related by

$$\frac{Q_1}{T_1} = S = \frac{Q_2}{T_2}. \quad (44.12)$$

But that means that if we have a single engine running between T_1 and T_2 , then the result of the whole analysis, the grand finale, is that Q_1 is to T_1 as Q_2 is to T_2 , if the engine absorbs energy Q_1 at temperature T_1 and delivers heat Q_2 at temperature T_2 . Whenever the engine is reversible, this relationship between the heats must follow. That is all there is to it: that is the center of the universe of thermodynamics.

If this is all there is to thermodynamics, why is it considered such a difficult subject? In doing a problem involving a given mass of some substance, the condition of the substance at any moment can be described by telling what its temperature is and what its volume is. If we know the temperature and volume of a substance, and that the pressure is some function of the temperature and volume, then we know the internal energy. One could say, "I do not want to do it that way. Tell me the temperature and the pressure, and I will tell you the volume. I can think of the volume as a function of temperature and pressure, and the internal energy as a function of temperature and pressure, and so on." That is why thermodynamics is hard, because everyone uses a different approach. If we could only sit down once and decide on our variables, and stick to them, it would be fairly easy.

Now we start to make deductions. Just as $F \Rightarrow ma$ is the center of the universe in mechanics, and it goes on and on and on after that, in the same way the principle just found is all there is to thermodynamics. But can one make conclusions out of it?

* We have previously defined our scale of temperature in a different way, namely by stating that the mean kinetic energy of a molecule in a perfect gas is proportional to the temperature, or that the perfect gas law says pV is proportional to T . Is this new definition equivalent? Yes, since the final result (44.7) derived from the gas law is the same as that derived here. We shall discuss this point again in the next chapter.

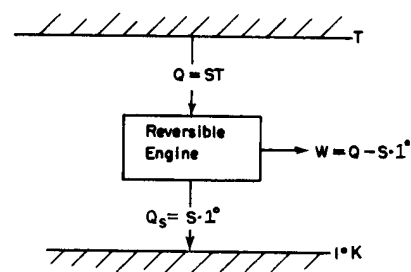


Fig. 44-9. Absolute thermodynamic temperature.

We begin. To obtain our first conclusion, we shall combine both laws, the law of conservation of energy and this law which relates the heats Q_2 and Q_1 , and we can easily obtain *the efficiency of a reversible engine*. From the first law, we have $W = Q_1 - Q_2$. According to our new principle,

$$Q_2 = \frac{T_2}{T_1} Q_1,$$

so the work becomes

$$W = Q_1 \left(1 - \frac{T_2}{T_1} \right) = Q_1 \frac{T_1 - T_2}{T_1}, \quad (44.13)$$

which tells us the efficiency of the engine—how much work we get out of so much heat. The efficiency of an engine is proportional to the difference in the temperatures between which the engine runs, divided by the higher temperature:

$$\text{Efficiency} = \frac{W}{Q_1} = \frac{T_1 - T_2}{T_1}. \quad (44.14)$$

The efficiency cannot be greater than unity and the absolute temperature cannot be less than zero, absolute zero. So, since T_2 must be positive, the efficiency is always less than unity. That is our first conclusion.

44-6 Entropy

Equation (44.7) or (44.12) can be interpreted in a special way. Working always with reversible engines, a heat Q_1 at temperature T_1 is “equivalent” to Q_2 at T_2 if $Q_1/T_1 = Q_2/T_2$, in the sense that as one is absorbed the other is delivered. This suggests that if we call Q/T something, we can say: in a reversible process as much Q/T is absorbed as is liberated; there is no gain or loss of Q/T . This Q/T is called *entropy*, and we say “there is no net change in entropy in a reversible cycle.” If T is 1° , then the entropy is $Q/1^\circ$ or, as we symbolized it, $Q_S/1^\circ = S$. Actually, S is the letter usually used for entropy, and it is numerically equal to the heat (which we have called Q_S) delivered to a 1° -reservoir (entropy is not itself a heat, it is heat divided by a temperature, hence it is measured in *joules per degree*).

Now it is interesting that besides the pressure, which is a function of the temperature and the volume, and the internal energy, which is a function of temperature and volume, we have found another quantity which is a function of the condition, i.e., the entropy of the substance. Let us try to explain how we compute it, and what we mean when we call it a “function of the condition.” Consider the system in two different conditions, much as we had in the experiment where we did the adiabatic and isothermal expansions. (Incidentally, there is no need that a heat engine have only two reservoirs, it could have three or four different temperatures at which it takes in and delivers heats, and so on.) We can move around on a pV diagram all over the place, and go from one condition to another. In other words, we could say the gas is in a certain condition a , and then it goes over to some other condition, b , and we will require that this transition, made from a to b , be reversible. Now suppose that all along the path from a to b we have little reservoirs at different temperatures, so that the heat dQ removed from the substance at each little step is delivered to each reservoir at the temperature corresponding to that point on the path. Then let us connect all these reservoirs, by reversible heat engines, to a single reservoir at the unit temperature. When we are finished carrying the substance from a to b , we shall bring all the reservoirs back to their original condition. Any heat dQ that has been absorbed from the substance at temperature T has now been converted by a reversible machine, and a certain amount of entropy dS has been delivered at the unit temperature as follows:

$$dS = dQ/T. \quad (44.15)$$

Let us compute the total amount of entropy which has been delivered. The entropy difference, or the entropy needed to go from a to b by this particular

reversible transformation, is the total entropy, the total of the entropy taken out of the little reservoirs, and delivered at the unit temperature:

$$S_b - S_a = \int_a^b \frac{dQ}{T}. \quad (44.16)$$

The question is, does the entropy difference depend upon the path taken? There is more than one way to go from a to b . Remember that in the Carnot cycle we could go from a to c in Fig. 44-6 by first expanding isothermally and then adiabatically; or we could first expand adiabatically and then isothermally. So the question is whether the entropy change which occurs when we go from a to b in Fig. 44-10 is the same on one route as it is on another. *It must be the same*, because if we went all the way around the cycle, going forward on one path and backward on another, we would have a reversible engine, and there would be no loss of heat to the reservoir at unit temperature. In a totally reversible cycle, no heat must be taken from the reservoir at the unit temperature, so the entropy needed to go from a to b is the same over one path as it is over another. It is *independent of path*, and depends only on the endpoints. We can, therefore, say that there is a certain function, which we call the entropy of the substance, that depends only on the condition, i.e., only on the volume and temperature.

We can find a function $S(V, T)$ which has the property that if we compute the change in entropy, as the substance is moved along any reversible path, in terms of the heat rejected at unit temperature, then

$$\Delta S = \int \frac{dQ}{T}, \quad (44.17)$$

where dQ is the heat removed from the substance at temperature T . This total entropy change is the difference between the entropy calculated at the initial and final points:

$$\Delta S = S(V_b, T_b) - S(V_a, T_a) = \int_a^b \frac{dQ}{T}. \quad (44.18)$$

This expression does not completely define the entropy, but rather only the *difference* of entropy between two different conditions. Only if we can evaluate the entropy for one special condition can we really define S absolutely.

For a long time it was believed that absolute entropy meant nothing—that only differences could be defined—but finally Nernst proposed what he called the *heat theorem*, which is also called the third law of thermodynamics. It is very simple. We will say what it is, but we will not explain why it is true. Nernst's postulate states simply that the entropy of any object at absolute zero is zero. We know of one case of T and V , namely $T = 0$, where S is zero; and so we can get the entropy at any other point.

To give an illustration of these ideas, let us calculate the entropy of a perfect gas. In an isothermal (and therefore reversible) expansion, $\int dQ/T$ is Q/T , since T is constant. Therefore (from 44.4) the change in entropy is

$$S(V_a, T) - S(V_b, T) = Nk \ln \frac{V_a}{V_b},$$

so $S(V, T) = Nk \ln V$ plus some function of T only. How does S depend on T ? We know that for a reversible adiabatic expansion, *no heat is exchanged*. Thus the entropy does not change even though V changes, provided that T changes also, such that $TV^{\gamma-1} = \text{constant}$. Can you see that this implies that

$$S(V, T) = Nk \left[\ln V + \frac{1}{\gamma - 1} \ln T \right] + a,$$

where a is some constant independent of *both* V and T ? [a is called the chemical constant. It depends on the gas in question, and may be determined experimentally from the Nernst theorem by measuring the heat liberated in cooling and condensing

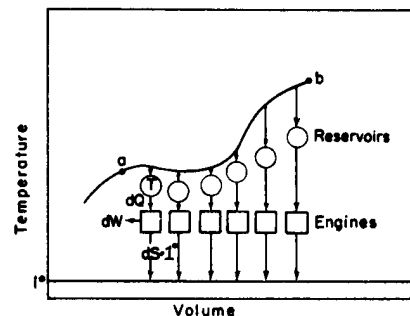


Fig. 44-10. Change in entropy during a reversible transformation.

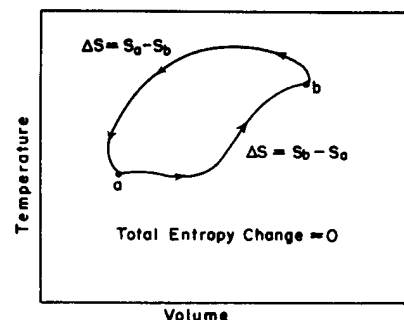


Fig. 44-11. Change in entropy in a completely reversible cycle.

the gas until it is brought to a solid (or for helium, a liquid) at 0° , by integrating $\int dQ/T$. It can also be determined theoretically by means of Planck's constant and quantum mechanics, but we shall not study it in this course.]

Now we shall remark on some of the properties of the entropy of things. We first remember that if we go along a reversible cycle from a to b , then the entropy of the substance will change by $S_b - S_a$. And we remember that as we go along the path, the entropy—the heat delivered at unit temperature—increases according to the rule $dS = dQ/T$, where dQ is the heat we remove from the substance when its temperature is T .

We already know that if we have a reversible cycle, the total entropy of everything is not changed, because the heat Q_1 absorbed at T_1 and the heat Q_2 delivered at T_2 correspond to equal and opposite changes in entropy, so that the net change in the entropy is zero. So for a reversible cycle there is no change in the entropy of anything, including the reservoirs. This rule may look like the conservation of energy again, but it is not; it applies only to reversible cycles. If we include irreversible cycles there is no law of conservation of entropy.

We shall give two examples. First, suppose that we do irreversible work on an object by friction, generating a heat Q on some object at temperature T . The entropy is increased by Q/T . The heat Q is equal to the work, and thus when we do a certain amount of work by friction against an object whose temperature is T , the entropy of the whole world increases by W/T .

Another example of irreversibility is this: If we put together two objects that are at different temperatures, say T_1 and T_2 , a certain amount of heat will flow from one to the other by itself. Suppose, for instance, we put a hot stone in cold water. Then when a certain heat ΔQ is transferred from T_1 to T_2 , how much does the entropy of the hot stone change? It decreases by $\Delta Q/T_1$. How much does the water entropy change? It increases by $\Delta Q/T_2$. The heat will, of course, flow only from the higher temperature T_1 to the lower temperature T_2 , so that ΔQ is positive if T_1 is greater than T_2 . So the change in entropy of the whole world is positive, and it is the difference of the two fractions:

$$\Delta S = \frac{\Delta Q}{T_2} - \frac{\Delta Q}{T_1}. \quad (44.19)$$

So the following proposition is true: in any process that is irreversible, the entropy of the whole world is increased. Only in reversible processes does the entropy remain constant. Since no process is absolutely reversible, there is always at least a small gain in the entropy; a reversible process is an idealization in which we have made the gain of entropy minimal.

Unfortunately, we are not going to enter into the field of thermodynamics very far. Our purpose is only to illustrate the principal ideas involved and the reasons why it is possible to make such arguments, but we will not use thermodynamics very much in this course. Thermodynamics is used very often by engineers and, particularly, by chemists. So we must learn our thermodynamics in practice in chemistry or engineering. Because it is not worth while duplicating everything, we shall just give some discussion of the origin of the theory, rather than much detail for special applications.

The two laws of thermodynamics are often stated this way:

First law: the energy of the universe is always constant.

Second law: the entropy of the universe is always increasing.

That is not a very good statement of the second law; it does not say, for example, that in a reversible cycle the entropy stays the same, and it does not say exactly what the entropy is. It is just a clever way of remembering the two laws, but it does not really tell us exactly where we stand. We have summarized the laws discussed in this chapter in Table 44-1. In the next chapter we shall apply these laws to discover the relationship between the heat generated in the expansion of a rubber band, and the extra tension when it is heated.

Table 44-1

Summary of the laws of thermodynamics

First law:

Heat put into a system + Work done on a system = Increase in internal energy of the system:

$$dQ + dW = dU.$$

Second law:

A process whose *only* net result is to take heat from a reservoir and convert it to work is impossible.

No heat engine taking heat Q_1 from T_1 and delivering heat Q_2 at T_2 can do more work than a reversible engine, for which

$$W = Q_1 - Q_2 = Q_1 \left(\frac{T_1 - T_2}{T_1} \right).$$

The entropy of a system is defined this way:

(a) If heat ΔQ is added reversibly to a system at temperature T , the increase in entropy of the system is $\Delta S = \Delta Q/T$.

(b) At $T = 0$, $S = 0$ (*third law*).

In a *reversible change*, the total entropy of all parts of the system (including reservoirs) does not change.

In *irreversible change*, the total entropy of the system always increases.

Illustrations of Thermodynamics

45-1 Internal energy

Thermodynamics is a rather difficult and complex subject when we come to apply it, and it is not appropriate for us to go very far into the applications in this course. The subject is of very great importance, of course, to engineers and chemists, and those who are interested in the subject can learn about the applications in physical chemistry or in engineering thermodynamics. There are also good reference books, such as Zemansky's *Heat and Thermodynamics*, where one can learn more about the subject. In the Encyclopedia Britannica, fourteenth edition, one can find excellent articles on thermodynamics and thermochemistry, and in the article on chemistry, the sections on physical chemistry, vaporization, liquefaction of gases, and so on.

The subject of thermodynamics is complicated because there are so many different ways of describing the same thing. If we wish to describe the behavior of a gas, we can say that the pressure depends on the temperature and on the volume, or we can say that the volume depends on the temperature and the pressure. Or with respect to the internal energy U , we might say that it depends on the temperature and volume, if those are the variables we have chosen—but we might also say that it depends on the temperature and the pressure, or the pressure and the volume, and so on. In the last chapter we discussed another function of temperature and volume, called the entropy S , and we can of course construct as many other functions of these variables as we like: $U - TS$ is a function of temperature and volume. So we have a large number of different quantities which can be functions of many different combinations of variables.

To keep the subject simple in this chapter, we shall decide at the start to use *temperature* and *volume* as the independent variables. Chemists use temperature and pressure, because they are easier to measure and control in chemical experiments, but we shall use temperature and volume throughout this chapter, except in one place where we shall see how to make the transformation into the chemists' system of variables.

We shall first, then, consider only one system of independent variables: temperature and volume. Secondly, we shall discuss only two dependent functions: the internal energy and the pressure. All the other functions can be derived from these, so it is not necessary to discuss them. With these limitations, thermodynamics is still a fairly difficult subject, but it is not quite so impossible!

First we shall review some mathematics. If a quantity is a function of two variables, the idea of the derivative of the quantity requires a little more careful thought than for the case where there is only one variable. What do we mean by the derivative of the pressure with respect to the temperature? The pressure change accompanying a change in the temperature depends partly, of course, on what happens to the *volume* while T is changing. We must specify the change in V before the concept of a derivative with respect to T has a precise meaning. We might ask, for example, for the rate of change of P with respect to T if V is held constant. This ratio is just the ordinary derivative that we usually write as dP/dT . We customarily use a special symbol, $\partial P/\partial T$, to remind us that P depends on another variable V as well as on T , and that this other variable is held constant. We shall not only use the symbol ∂ to call attention to the fact that the other variable is held constant, but we shall also write the variable that is held constant as a subscript, $(\partial P/\partial T)_V$. Since we have only two independent variables, this notation is redundant, but it will help us keep our wits about us in the thermodynamic jungle of partial derivatives.

45-1 Internal energy

45-2 Applications

45-3 The Clausius-Clapeyron equation

Let us suppose that the function $f(x, y)$ depends on the two independent variables x and y . By $(\partial f/\partial x)_y$ we mean simply the ordinary derivative, obtained in the usual way, if we treat y as a constant:

$$\left(\frac{\partial f}{\partial x}\right)_y = \lim_{\Delta x \rightarrow 0} \frac{f(x + \Delta x, y) - f(x, y)}{\Delta x}.$$

Similarly, we define

$$\left(\frac{\partial f}{\partial y}\right)_x = \lim_{\Delta y \rightarrow 0} \frac{f(x, y + \Delta y) - f(x, y)}{\Delta y}.$$

For example, if $f(x, y) = x^2 + yx$, then $(\partial f/\partial x)_y = 2x + y$, and $(\partial f/\partial y)_x = x$. We can extend this idea to higher derivatives: $\partial^2 f/\partial y^2$ or $\partial^2 f/\partial y \partial x$. The latter symbol indicates that we first differentiate f with respect to x , treating y as a constant, then differentiate the result with respect to y , treating x as a constant. The actual order of differentiation is immaterial: $\partial^2 f/\partial x \partial y = \partial^2 f/\partial y \partial x$.

We will need to compute the change Δf in $f(x, y)$ when x changes to $x + \Delta x$ and y changes to $y + \Delta y$. We assume throughout the following that Δx and Δy are infinitesimally small:

$$\begin{aligned} \Delta f &= f(x + \Delta x, y + \Delta y) - f(x, y) \\ &= \underbrace{f(x + \Delta x, y + \Delta y) - f(x, y + \Delta y)}_{\Delta x \left(\frac{\partial f}{\partial x}\right)_y} + \underbrace{f(x, y + \Delta y) - f(x, y)}_{\Delta y \left(\frac{\partial f}{\partial y}\right)_x} \\ &= \Delta x \left(\frac{\partial f}{\partial x}\right)_y + \Delta y \left(\frac{\partial f}{\partial y}\right)_x. \end{aligned} \quad (45.1)$$

The last equation is the fundamental relation that expresses Δf in terms of Δx and Δy .

As an example of the use of this relation, let us calculate the change in the internal energy $U(T, V)$ when the temperature changes from T to $T + \Delta T$ and the volume changes from V to $V + \Delta V$. Using Eq. (45.1), we write

$$\Delta U = \Delta T \left(\frac{\partial U}{\partial T}\right)_V + \Delta V \left(\frac{\partial U}{\partial V}\right)_T. \quad (45.2)$$

In our last chapter we found another expression for the change ΔU in the internal energy when a quantity of heat ΔQ was added to the gas:

$$\Delta U = \Delta Q - P \Delta V. \quad (45.3)$$

In comparing Eqs. (45.2) and (45.3) one might at first be inclined to think that $P = (\partial U/\partial V)_T$, but this is not correct. To obtain the correct relation, let us first suppose that we add a quantity of heat ΔQ to the gas while keeping the volume constant, so that $\Delta V = 0$. With $\Delta V = 0$, Eq. (45.3) tells us that $\Delta U = \Delta Q$, and Eq. (45.2) tells us that $\Delta U = (\partial U/\partial T)_V \Delta T$, so that $(\partial U/\partial T)_V = \Delta Q/\Delta T$. The ratio $\Delta Q/\Delta T$, the amount of heat one must put into a substance in order to change its temperature by one degree with the volume held constant, is called the *specific heat at constant volume* and is designated by the symbol C_V . By this argument we have shown that

$$\left(\frac{\partial U}{\partial T}\right)_V = C_V. \quad (45.4)$$

Now let us again add a quantity of heat ΔQ to the gas, but this time we will hold T constant and allow the volume to change by ΔV . The analysis in this case is more complex, but we can calculate ΔU by the argument of Carnot, making use of the Carnot cycle we introduced in the last chapter.

The pressure-volume diagram for the Carnot cycle is shown in Fig. 45-1. As we have already shown, the total amount of work done by the gas in a reversible cycle is $\Delta Q(\Delta T/T)$, where ΔQ is the amount of heat energy added to the gas as it expands isothermally at temperature T from volume V to $V + \Delta V$, and $T - \Delta T$ is the final temperature reached by the gas as it expands adiabatically on the second leg of the cycle. Now we will show that this work done is also given by

the shaded area in Fig. 45-1. In any circumstances, the work done by the gas is $\int P dV$, and is positive when the gas expands and negative when the gas is compressed. If we plot P vs. V , the variation of P and V is represented by a curve which gives the value of P corresponding to a particular value of V . As the volume changes from one value to another, the work done by the gas, the integral $\int P dV$, is the area under the curve connecting the initial and final values of V . When we apply this idea to the Carnot cycle, we see that as we go around the cycle, paying attention to the sign of the work done by the gas, the net work done by the gas is just the shaded area in Fig. 45-1.

Now we want to evaluate the shaded area geometrically. The cycle we have used in Fig. 45-1 differs from that used in the previous chapter in that we now suppose that ΔT and ΔQ are infinitesimally small. We are working between adiabatic lines and isothermal lines that are very close together, and the figure described by the heavy lines in Fig. 45-1 will approach a parallelogram as the increments ΔT and ΔQ approach zero. The area of this parallelogram is just $\Delta V \Delta P$, where ΔV is the change in volume as energy ΔQ is added to the gas at constant temperature, and ΔP is the change in pressure as the temperature changes by ΔT at constant volume. One can easily show that the shaded area in Fig. 45-1 is given by $\Delta V \Delta P$ by recognizing that the shaded area is equal to the area enclosed by the dotted lines in Fig. 45-2, which in turn differs from the rectangle bounded by ΔP and ΔV only by the addition and subtraction of the equal triangular areas in Fig. 45-2.

Now let us summarize the results of the arguments we have developed so far:

$$\left. \begin{aligned} &\text{Work done by the gas} = \text{shaded area} = \Delta V \Delta P = \Delta Q \left(\frac{\Delta T}{T} \right) \\ \text{or} & \quad \frac{\Delta T}{T} \cdot (\text{heat needed to change } V \text{ by } \Delta V)_{\text{constant } T} \\ &= \Delta V \cdot (\text{change in } P \text{ when } T \text{ changes by } \Delta T)_{\text{constant } V} \\ \text{or} & \quad \frac{1}{\Delta V} \cdot (\text{heat needed to change } V \text{ by } \Delta V)_T = T \left(\frac{\partial P}{\partial T} \right)_V. \end{aligned} \right\} \quad (45.5)$$

Equation (45.5) expresses the essential result of Carnot's argument. The whole of thermodynamics can be deduced from Eq. (45.5) and the First Law, which is stated in Eq. (45.3). Equation (45.5) is essentially the Second Law, although it was originally deduced by Carnot in a slightly different form, since he did not use our definition of temperature.

Now we can proceed to calculate $(\partial U / \partial V)_T$. By how much would the internal energy U change if we changed the volume by ΔV ? First, U changes because heat is put in, and second, U changes because work is done. The heat put in is

$$\Delta Q = T \left(\frac{\partial P}{\partial T} \right)_V \Delta V,$$

according to Eq. (45.5), and the work done on the substance is $-P \Delta V$. Therefore the change ΔU in internal energy has two pieces:

$$\Delta U = T \left(\frac{\partial P}{\partial T} \right)_V \Delta V - P \Delta V. \quad (45.6)$$

Dividing both sides by ΔV , we find for the rate of change of U with V at constant T

$$\left(\frac{\partial U}{\partial V} \right)_T = T \left(\frac{\partial P}{\partial T} \right)_V - P. \quad (45.7)$$

In our thermodynamics, in which T and V are the only variables and P and U are the only functions, Eqs. (45.3) and (45.7) are the basic equations from which all the result of the subject can be deduced.

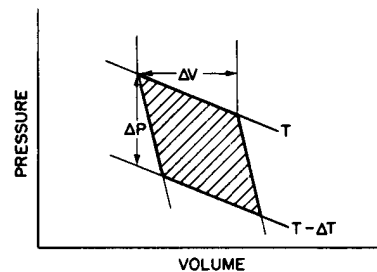


Fig. 45-1. Pressure-volume diagram for a Carnot cycle. The curves marked T and $T - \Delta T$ are isothermal lines; the steeper curves are adiabatic lines. ΔV is the volume change as heat ΔQ is added to the gas at constant temperature T . ΔP is the pressure change at constant volume as the gas temperature is changed from T to $T - \Delta T$.

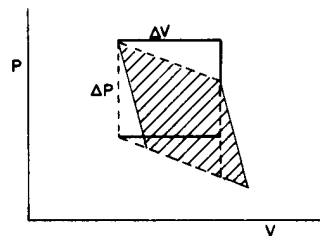


Fig. 45-2. Shaded area = area enclosed by dashed lines = area of rectangle = $\Delta P \Delta V$.

45-2 Applications

Now let us discuss the meaning of Eq. (45.7) and see why it answers the questions which we proposed in our last chapter. We considered the following problem: in kinetic theory it is obvious that an increase in temperature leads to an increase in pressure, because of the bombardments of the atoms on a piston. For the same physical reason, when we let the piston move back, heat is taken out of the gas and, in order to keep the temperature constant, heat will have to be put back in. The gas cools when it expands, and the pressure rises when it is heated. There must be some connection between these two phenomena, and this connection is given explicitly in Eq. (45.7). If we hold the volume fixed and increase the temperature, the pressure rises at a rate $(\partial P/\partial T)_V$. Related to that fact is this: if we increase the volume, the gas will cool unless we pour some heat in to maintain the temperature constant, and $(\partial U/\partial V)_T$ tells us the amount of heat needed to maintain the temperature. Equation (45.7) expresses the fundamental interrelationship between these two effects. That is what we promised we would find when we came to the laws of thermodynamics. Without knowing the internal mechanism of the gas, and knowing only that we cannot make perpetual motion of the second type, we can deduce the relationship between the amount of heat needed to maintain a constant temperature when the gas expands, and the pressure change when the gas is heated!

Now that we have the result we wanted for a gas, let us consider the rubber band. When we stretch a rubber band, we find that its temperature falls, and when we heat a rubber band, we find that it pulls itself in. What is the equation that gives the same relation for a rubber band as Eq. (45.3) gives for gas? For a rubber band the situation will be something like this: when heat ΔQ is put in, the internal energy is changed by ΔU and some work is done. The only difference will be that the work done by the rubber band is $-F \Delta L$ instead of $P \Delta V$, where F is the force on the band, and L is the length of the band. The force F is a function of temperature and of length of the band. Replacing $P \Delta V$ in Eq. (45.3) by $-F \Delta L$, we get

$$\Delta U = \Delta Q + F \Delta L \quad (45.8)$$

Comparing Eqs. (45.3) and (45.8), we see that the rubber band equation is obtained by a mere substitution of one letter for another. Furthermore, if we substitute L for V , and $-F$ for P , all of our discussion of the Carnot cycle applies to the rubber band. We can immediately deduce, for instance, that the heat ΔQ needed to change the length by ΔL is given by the analog to Eq. (45.5): $\Delta Q = -T(\partial F/\partial T)_L \Delta L$. This equation tells us that if we keep the length of a rubber band fixed and heat the band, we can calculate how much the force will increase in terms of the heat needed to keep the temperature constant when the band is stretched a little bit. So we see that the same equation applies to both gas and a rubber band. In fact, if one can write $\Delta U = \Delta Q + A \Delta B$, where A and B represent different quantities, force and length, pressure and volume, etc., one can apply the results obtained for a gas by substituting A and B for P and V . For example, consider the electric potential difference, or "voltage," E in a battery and the charge ΔZ that moves through the battery. We know that the work done in a reversible electric cell, like a storage battery, is $E \Delta Z$. (Since we include no $P \Delta V$ term in the work, we require that our battery maintain a constant volume.) Let us see what thermodynamics can tell us about the performance of a battery. If we substitute E for P and Z for V in Eq. (45.6), we obtain

$$\frac{\Delta U}{\Delta Z} = -T \left(\frac{\partial E}{\partial T} \right)_Z - E. \quad (45.9)$$

Equation (45.9) says that the internal energy U is changed when a charge ΔZ moves through the cell. Why is $\Delta U/\Delta Z$ not simply the voltage E of the battery? The answer is that a real battery gets warm when charge moves through the cell. The internal energy of the battery is changed, first, because the battery did some work on the outside circuit, and second, because the battery is heated. The re-

markable thing is that the second part can again be expressed in terms of the way in which the battery voltage changes with temperature. Incidentally, when the charge moves through the cell, chemical reactions occur, and Eq. (45.9) suggests a nifty way of measuring the amount of energy required to produce a chemical reaction. All we need do is construct a cell that works on the reaction, measure the voltage, and measure how much the voltage changes with temperature when we draw no charge from the battery!

Now we have assumed that the volume of the battery can be maintained constant, since we have omitted the $P\Delta V$ term when we set the work done by the battery equal to $E\Delta Z$. It turns out that it is technically quite difficult to keep the volume constant. It is much easier to keep the cell at constant atmospheric pressure. For that reason, the chemists do not like any of the equations we have written above: they prefer equations which describe performance under constant pressure. We chose at the beginning of this chapter to use V and T as independent variables. The chemists prefer P and T , and we will now consider how the results we have obtained so far can be transformed into the chemists' system of variables. Remember that in the following treatment confusion can easily set in because we are shifting gears from T and V to T and P .

We started in Eq. (45.3) with $\Delta U = \Delta Q - P\Delta V$; $P\Delta V$ may be replaced by $E\Delta Z$ or $A\Delta B$. If we could somehow replace the last term, $P\Delta V$, by $V\Delta P$, then we would have interchanged V and P , and the chemists would be happy. Well, a clever man noticed that the differential of the product PV is $d(PV) = PdV + VdP$, and if he added this equality to Eq. (45.3), he obtained

$$\begin{array}{rcl} \Delta(PV) & = & P\Delta V + V\Delta P \\ \Delta U & = & \Delta Q - P\Delta V \\ \hline \Delta(U + PV) & = & \Delta Q + V\Delta P \end{array}$$

In order that the result look like Eq. (45.3), we define $U + PV$ to be something new, called the *enthalpy*, H , and we write $\Delta H = \Delta Q + V\Delta P$.

Now we are ready to transform our results into chemists' language with the following rules: $U \rightarrow H$, $P \rightarrow -V$, $V \rightarrow P$. For example, the fundamental relationship that chemists would use instead of Eq. (45.7) is

$$\left(\frac{\partial H}{\partial P}\right)_T = T\left(\frac{\partial V}{\partial T}\right)_P - V.$$

It should now be clear how one transforms to the chemists' variables T and P . We now go back to our original variables: for the remainder of this chapter, T and V are the independent variables.

Now let us apply the results we have obtained to a number of physical situations. Consider first the ideal gas. From kinetic theory we know that the internal energy of a gas depends only on the motion of the molecules and the number of molecules. The internal energy depends on T , but not on V . If we change V , but keep T constant, U is not changed. Therefore $(\partial U/\partial V)_T = 0$, and Eq. (45.7) tells us that for an ideal gas

$$T\left(\frac{\partial P}{\partial T}\right)_V - P = 0. \quad (45.10)$$

Equation (45.10) is a differential equation that can tell us something about P . We take account of the partial derivatives in the following way: Since the partial derivative is at constant V , we will replace the partial derivative by an ordinary derivative and write explicitly, to remind us, "constant V ." Equation (45.10) then becomes

$$T \frac{\Delta P}{\Delta T} - P = 0; \quad \text{const } V, \quad (45.11)$$

which we can integrate to get

$$\begin{array}{ll} \ln P = \ln T + \text{const}; & \text{const } V, \\ P = \text{const} \times T; & \text{const } V. \end{array} \quad (45.12)$$

We know that for an ideal gas the pressure is equal to

$$P = \frac{RT}{V}, \quad (45.13)$$

which is consistent with (45.12), since V and R are constants. Why did we bother to go through this calculation if we already knew the results? Because we have been using *two independent definitions of temperature!* At one stage we assumed that the kinetic energy of the molecules was proportional to the temperature, an assumption that defines one scale of temperature which we will call the ideal gas scale. The T in Eq. (45.13) is based on the gas scale. We also call temperatures measured on the gas scale *kinetic* temperatures. Later, we defined the temperature in a second way which was completely independent of any substance. From arguments based on the Second Law we defined what we might call the “grand thermodynamic absolute temperature” T , the T that appears in Eq. (45.12). What we proved here is that the pressure of an ideal gas (defined as one for which the internal energy does not depend on the volume) is proportional to the grand thermodynamic absolute temperature. We also know that the pressure is proportional to the temperature measured on the gas scale. Therefore we can deduce that the kinetic temperature is proportional to the “grand thermodynamic absolute temperature.” That means, of course, that if we were sensible we could make two scales agree. In this instance, at least, the two scales *have* been chosen so that they coincide; the proportionality constant has been chosen to be 1. Most of the time man chooses trouble for himself, but in this case he made them equal!

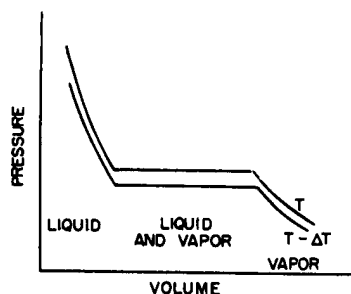


Fig. 45-3. Isothermal lines for a condensable vapor compressed in a cylinder. At the left, the substance is in the liquid phase. At the right, the substance is vaporized. In the center, both liquid and vapor are present in the cylinder.

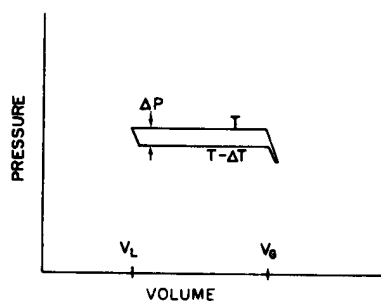


Fig. 45-4. Pressure-volume diagram for a Carnot cycle with a condensable vapor in the cylinder. At the left, the substance is in the liquid state. A quantity of heat L is added at temperature T to vaporize the liquid. The vapor expands adiabatically as T changes to $T - \Delta T$.

45-3 The Clausius-Clapeyron equation

The vaporization of a liquid is another application of the results we have derived. Suppose we have some liquid in a cylinder, such that we can compress it by pushing on the piston, and we ask ourselves, “If we keep the temperature constant, how does the pressure vary with volume?” In other words, we want to draw an isothermal line on the P - V diagram. The substance in the cylinder is not the ideal gas that we considered earlier; now it may be in the liquid or the vapor phase, or both may be present. If we apply sufficient pressure, the substance will condense to a liquid. Now if we squeeze still harder, the volume changes very little, and our isothermal line rises rapidly with decreasing volume, as shown at the left in Fig. 45-3.

If we increase the volume by pulling the piston out, the pressure drops until we reach the point at which the liquid starts to boil, and then vapor starts to form. If we pull the piston out farther, all that happens is that more liquid vaporizes. When there is part liquid and part vapor in the cylinder, the two phases are in equilibrium—liquid is evaporating and vapor is condensing at the same rate. If we make more room for the vapor, more vapor is needed to maintain the pressure, so a little more liquid evaporates, but the pressure remains constant. On the flat part of the curve in Fig. 45-3 the pressure does not change, and the value of the pressure here is called the *vapor pressure at temperature T* . As we continue to increase the volume, there comes a time when there is no more liquid to evaporate. At this juncture, if we expand the volume further, the pressure will fall as for an ordinary gas, as shown at the right of the P - V diagram. The lower curve in Fig. 45-3 is the isothermal line at a slightly lower temperature $T - \Delta T$. The pressure in the liquid phase is slightly reduced because liquid expands with an increase in temperature (for most substances, but not for water near the freezing point) and, of course, the vapor pressure is lower at the lower temperature.

We will now make a cycle out of the two isothermal lines by connecting them (say by adiabatic lines) at the ends of the flat sections, as shown in Fig. 45-4. The little jiggle in the lower right-hand corner of the figure will make little difference and we will neglect it. We are going to use the argument of Carnot, which tells us that the heat added to the substance in changing it from a liquid to a vapor is related to the work done by the substance as it goes around the cycle. Let

us call L the heat needed to vaporize the substance in the cylinder. As in the argument immediately preceding Eq. (45.5), we know that $L(\Delta T/T) = \text{work done by the substance}$. As before, the work done by the substance is the shaded area, which is approximately $\Delta P(V_G - V_L)$, where ΔP is the difference in vapor pressure at the two temperatures T and $T - \Delta T$, V_G is the volume of the gas, and V_L is the volume of the liquid, both volumes measured at the vapor pressure. Setting these two expressions for the area equal, we get $L\Delta T/T = \Delta P(V_G - V_L)$, or

$$\frac{L}{T(V_G - V_L)} = (\partial P_{\text{vap}}/\partial T). \quad (45.14)$$

Equation (45.14) gives the relationship between the rate of change of vapor pressure with temperature and the amount of heat required to evaporate the liquid. This relationship was deduced by Carnot, but it is called the Clausius-Clapeyron equation.

Now let us compare Eq. (45.14) with the results deduced from kinetic theory. Usually V_G is very much larger than V_L . So $V_G - V_L \approx V_G = RT/P$ per mole. If we further assume that L is a constant, independent of temperature—not a very good approximation—then we would have $\partial P/\partial T = L/(RT^2P)$. The solution of this differential equation is

$$P = \text{const } e^{-L/RT}. \quad (45.15)$$

Let us compare this with the pressure variation with temperature that we deduced earlier from kinetic theory. Kinetic theory indicated the possibility, at least roughly, that the number of molecules of vapor above a liquid would be

$$n = \left(\frac{1}{V_A}\right) e^{-(U_G - U_L)/RT}, \quad (45.16)$$

where $U_G - U_L$ is the internal energy per mole in the liquid minus the internal energy per mole in the gas, i.e., the energy needed to vaporize a mole of liquid. Equation (45.15) from thermodynamics and Eq. (45.16) from kinetic theory are very closely related because the pressure is nkT , but they are not exactly the same. However, they will turn out to be exactly the same if we assume $L - U_G = \text{const}$, instead of $L = \text{const}$. If we assume $L - U_G = \text{const}$, independent of temperature, then the argument leading to Eq. (45.15) will produce Eq. (45.16).

This comparison shows the advantages and disadvantages of thermodynamics over kinetic theory: First of all, Eq. (45.14) obtained by thermodynamics is exact, while Eq. (45.16) can only be approximated, for instance, if U is nearly constant, and if the model is right. Second, we may not understand correctly how the gas goes into the liquid; nevertheless, Eq. (45.14) is right, while (45.16) is only approximate. Third, although our treatment applies to a gas condensing into a liquid, the argument is true for any other change of state. For instance, the solid-to-liquid transition has the same kind of curve as that shown in Figs. 45-3 and 45-4. Introducing the latent heat for melting, M/mole , the formula analogous to Eq. (45.14) then is $(\partial P_{\text{melt}}/\partial T)_V = M/[T(V_{\text{liq}} - V_{\text{solid}})]$. Although we may not understand the kinetic theory of the melting process, we nevertheless have a correct equation. However, when we *can* understand the kinetic theory, we have another advantage. Equation (45.14) is only a differential relationship, and we have no way of obtaining the constants of integration. In the kinetic theory we can obtain the constants also if we have a good model that describes the phenomenon completely. So there are advantages and disadvantages to each. When knowledge is weak and the situation is complicated, thermodynamic relations are really the most powerful. When the situation is very simple and a theoretical analysis can be made, then it is better to try to get more information from theoretical analysis.

One more example: blackbody radiation. We have discussed a box containing radiation and nothing else. We have talked about the equilibrium between the oscillator and the radiation. We also found that the photons hitting the wall of the box would exert the pressure P , and we found $PV = U/3$, where U is the

total energy of all the photons and V is the volume of the box. If we substitute $U = 3PV$ in the basic Eq. (45.7), we find

$$\left(\frac{\partial U}{\partial V}\right)_T = 3P = T\left(\frac{\partial P}{\partial T}\right)_V - P. \quad (45.17)$$

Since the volume of our box is constant, we can replace $(\partial P/\partial T)_V$ by dP/dT to obtain an ordinary differential equation we can integrate: $\ln P = 4 \ln T + \text{const}$, or $P = \text{const} \times T^4$. The pressure of radiation varies as the fourth power of the temperature, and the energy content of the radiation, $U/V = P/3$, also varies as T^4 . It is usual to write $U/V = (4\sigma/c)T^4$, where c is the speed of light and σ is a constant. It is not possible to get σ from thermodynamics alone. Here is a good example of its power, and its limitations. To know that U/V goes as T^4 is a great deal, but to know how big U/V actually is at any temperature requires that we go into the kind of detail that only a complete theory can supply. For blackbody radiation we have such a theory and we can derive an expression for the constant σ in the following manner.

Let $I(\omega) d\omega$ be the intensity distribution, the energy flow through 1 m^2 in one second with frequency between ω and $\omega + d\omega$. The energy density distribution = energy/volume = $I(\omega) d\omega/c$ is

$$\begin{aligned} \frac{U}{V} &= \text{total energy density} \\ &= \int_{\omega=0}^{\infty} \text{energy density } \omega \text{ and } \omega + d\omega \\ &= \int_0^{\infty} \frac{I(\omega) d\omega}{c}. \end{aligned}$$

From our earlier discussions, we know that

$$I(\omega) = \frac{\hbar\omega^3}{\pi^2 c^2 (e^{\hbar\omega/kT} - 1)}.$$

Substituting this expression for $I(\omega)$ in our equation for U/V , we get

$$\frac{U}{V} = \frac{1}{\pi^2 c^3} \int_0^{\infty} \frac{\hbar\omega^3 d\omega}{e^{\hbar\omega/kT} - 1}.$$

If we substitute $x = \hbar\omega/kT$, the expression becomes

$$\frac{U}{V} = \frac{(kT)^4}{\hbar^3 \pi^2 c^3} \int_0^{\infty} \frac{x^3 dx}{e^x - 1}.$$

This integral is just some number that we can get, approximately, by drawing a curve and taking the area by counting squares. It is roughly 6.5. The mathematicians among us can show that the integral is exactly $\pi^4/15$.* Comparing this expression with $U/V = (4\sigma/c)T^4$, we find

$$\sigma = \frac{k^4 \pi^2}{60 \hbar^3 c^2} = 5.67 \times 10^{-8} \frac{\text{watts}}{(\text{meter})^2 (\text{degree})^4}.$$

* Since $(e^x - 1)^{-1} = e^{-x} + e^{-2x} + \dots$, the integral is

$$\sum_{n=1}^{\infty} \int_0^{\infty} e^{-nx} x^3 dx.$$

But $\int_0^{\infty} e^{-nx} dx = 1/n$, and differentiating with respect to n three times gives $\int_0^{\infty} x^3 e^{-nx} dx = 6/n^4$, so the integral is $6(1 + \frac{1}{16} + \frac{1}{81} + \dots)$ and a good estimate comes from adding the first few terms. In Chapter 50 we will find a way to show that the sum of the reciprocal fourth powers of the integers is, in fact, $\pi^4/90$.

If we make a small hole in our box, how much energy will flow per second through the hole of unit area? To go from energy density to energy flow, we multiply the energy density U/V by c . We also multiply by $\frac{1}{4}$, which arises as follows: first, a factor of $\frac{1}{2}$, because only the energy which is flowing *out* escapes; and second, another factor $\frac{1}{2}$, because energy which approaches the hole at an angle to the normal is less effective in getting through the hole by a cosine factor. The average value of the cosine is $\frac{1}{2}$. It is clear now why we write $U/V = (4\sigma/c)T^4$: so that we can ultimately say that the flux from a small hole is σT^4 per unit area.

Ratchet and pawl

46-1 How a ratchet works

In this chapter we discuss the ratchet and pawl, a very simple device which allows a shaft to turn only one way. The possibility of having something turn only one way requires some detailed and careful analysis, and there are some very interesting consequences.

The plan of the discussion came about in attempting to devise an elementary explanation, from the molecular or kinetic point of view, for the fact that there is a maximum amount of work which can be extracted from a heat engine. Of course we have seen the essence of Carnot's argument, but it would be nice to find an explanation which is elementary in the sense that we can see what is happening physically. Now, there are complicated mathematical demonstrations which follow from Newton's laws to demonstrate that we can get only a certain amount of work out when heat flows from one place to another, but there is great difficulty in converting this into an elementary demonstration. In short, we do not understand it, although we can follow the mathematics.

In Carnot's argument, the fact that more than a certain amount of work cannot be extracted in going from one temperature to another is deduced from another axiom, which is that if everything is at the same temperature, heat cannot be converted to work by means of a cyclic process. First, let us back up and try to see, in at least one elementary example, why this simpler statement is true.

Let us try to invent a device which will violate the Second Law of Thermodynamics, that is, a gadget which will generate work from a heat reservoir with everything at the same temperature. Let us say we have a box of gas at a certain temperature, and inside there is an axle with vanes in it. (See Fig. 46-1 but take $T_1 = T_2 = T$, say.) Because of the bombardments of gas molecules on the vane, the vane oscillates and jiggles. All we have to do is to hook onto the other end of the axle a wheel which can turn only one way—the ratchet and pawl. Then when the shaft tries to jiggle one way, it will not turn, and when it jiggles the other, it will turn. Then the wheel will slowly turn, and perhaps we might even tie a flea onto a string hanging from a drum on the shaft, and lift the flea! Now let us ask if this is possible. According to Carnot's hypothesis, it is impossible. But if we just look at it, we see, *prima facie*, that it seems quite possible. So we must look more closely. Indeed, if we look at the ratchet and pawl, we see a number of complications.

First, our idealized ratchet is as simple as possible, but even so, there is a pawl, and there must be a spring in the pawl. The pawl must return after coming off a tooth, so the spring is necessary.

Another feature of this ratchet and pawl, not shown in the figure, is quite essential. Suppose the device were made of perfectly elastic parts. After the pawl is lifted off the end of the tooth and is pushed back by the spring, it will bounce against the wheel and continue to bounce. Then, when another fluctuation came, the wheel could turn the other way, because the tooth could get underneath during the moment when the pawl was up! Therefore an essential part of the irreversibility of our wheel is a damping or deadening mechanism which stops the bouncing. When the damping happens, of course, the energy that was in the pawl goes into the wheel and shows up as heat. So, as it turns, the wheel will get hotter and hotter. To make the thing simpler, we can put a gas around the wheel to take up some of the heat. Anyway, let us say the gas keeps rising in temperature, along with the wheel. Will it go on forever? No! The pawl and wheel, both at some temperature

46-1 How a ratchet works

46-2 The ratchet as an engine

46-3 Reversibility in mechanics

46-4 Irreversibility

46-5 Order and entropy

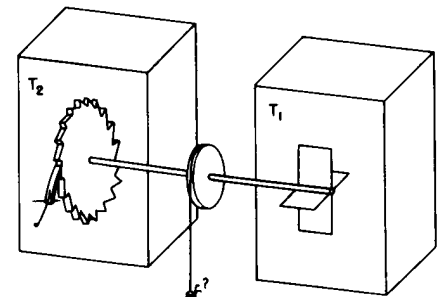


Fig. 46-1. The ratchet and pawl machine.

T , also have Brownian motion. This motion is such that, every once in a while, by accident, the pawl lifts itself up and over a tooth just at the moment when the Brownian motion on the vanes is trying to turn the axle backwards. And as things get hotter, this happens more often.

So, this is the reason this device does not work in perpetual motion. When the vanes get kicked, sometimes the pawl lifts up and goes over the end. But sometimes, when it tries to turn the other way, the pawl has already lifted due to the fluctuations of the motions on the wheel side, and the wheel goes back the other way! The net result is nothing. It is not hard to demonstrate that when the temperature on both sides is equal, there will be no net average motion of the wheel. Of course the wheel will do a lot of jiggling this way and that way, but it will not do what we would like, which is to turn just one way.

Let us look at the reason. It is necessary to do work against the spring in order to lift the pawl to the top of a tooth. Let us call this energy ϵ , and let θ be the angle between the teeth. The chance that the system can accumulate enough energy, ϵ , to get the pawl over the top of the tooth, is $e^{-\epsilon/kT}$. But the probability that the pawl will accidentally be up is also $e^{-\epsilon/kT}$. So the number of times that the pawl is up and the wheel can turn backwards freely is equal to the number of times that we have enough energy to turn it forward when the pawl is down. We thus get a "balance," and the wheel will not go around.

46-2 The ratchet as an engine

Let us now go further. Take the example where the temperature of the vanes is T_1 and the temperature of the wheel, or ratchet, is T_2 , and T_2 is less than T_1 . Because the wheel is cold and the fluctuations of the pawl are relatively infrequent, it will be very hard for the pawl to attain an energy ϵ . Because of the high temperature T_1 , the vanes will often attain the energy ϵ , so our gadget will go in one direction, as designed.

We would now like to see if it can lift weights. Onto the drum in the middle we tie a string, and put a weight, such as our flea, on the string. We let L be the torque due to the weight. If L is not too great, our machine will lift the weight because the Brownian fluctuations make it more likely to move in one direction than the other. We want to find how much weight it can lift, how fast it goes around, and so on.

First we consider a forward motion, the usual way one designs a ratchet to run. In order to make one step forward, how much energy has to be borrowed from the vane end? We must borrow an energy ϵ to lift the pawl. The wheel turns through an angle θ against a torque L , so we also need the energy $L\theta$. The total amount of energy that we have to borrow is thus $\epsilon + L\theta$. The probability that we get this energy is proportional to $e^{-(\epsilon+L\theta)/kT_1}$. Actually, it is not only a question of getting the energy, but we also would like to know the number of times per second it has this energy. The probability per second is proportional to $e^{-(\epsilon+L\theta)/kT_1}$, and we shall call the proportionality constant $1/\tau$. It will cancel out in the end anyway. When a forward step happens, the work done on the weight is $L\theta$. The energy taken from the vane is $\epsilon + L\theta$. The spring gets wound up with energy ϵ , then it goes clatter, clatter, bang, and this energy goes into heat. All the energy taken out goes to lift the weight and to drive the pawl, which then falls back and gives heat to the other side.

Now we look at the opposite case, which is backward motion. What happens here? To get the wheel to go backwards all we have to do is supply the energy to lift the pawl high enough so that the ratchet will slip. This is still energy ϵ . Our probability per second for the pawl to lift this high is now $(1/\tau)e^{-\epsilon/kT_2}$. Our proportionality constant is the same, but this time kT_2 shows up because of the different temperature. When this happens, the work is released because the wheel slips backward. It loses one notch, so it releases work $L\theta$. The energy taken from the ratchet system is ϵ , and the energy given to the gas at T_1 on the vane side is $L\theta + \epsilon$. It takes a little thinking to see the reason for that. Suppose the pawl has lifted itself up accidentally by a fluctuation. Then when it falls back and the spring

Table 46-1

Summary of operation of ratchet and pawl.

<i>Forward:</i>	Need energy	$\epsilon + L\theta$	from vane. \therefore Rate = $\frac{1}{\tau} e^{-(L\theta + \epsilon)/kT_1}$
	Takes from vane	$L\theta + \epsilon$	
	Does work	$L\theta$	
	Gives to ratchet	ϵ	

<i>Backward:</i>	Needs energy	ϵ	for pawl. \therefore Rate = $\frac{1}{\tau} e^{-\epsilon/kT_2}$
	Takes from ratchet	ϵ	} same as above with sign reversed.
	Releases work	$L\theta$	
	Gives to vane	$L\theta + \epsilon$	

If system is reversible, rates are equal, hence $\frac{\epsilon + L\theta}{T_1} = \frac{\epsilon}{T_2}$.

$\frac{\text{Heat to ratchet}}{\text{Heat from vane}} = \frac{\epsilon}{L\theta + \epsilon}$. Hence $\frac{Q_2}{Q_1} = \frac{T_2}{T_1}$.

pushes it down against the tooth, there is a force trying to turn the wheel, because the tooth is pushing on an inclined plane. This force is doing work, and so is the force due to the weights. So both together make up the total force, and all the energy which is slowly released appears at the vane end as heat. (Of course it must, by conservation of energy, but one must be careful to think the thing through!) We notice that all these energies are exactly the same, but reversed. So, depending upon which of these two rates is greater, the weight is either slowly lifted or slowly released. Of course, it is constantly jiggling around, going up for a while and down for a while, but we are talking about the average behavior.

Suppose that for a particular weight the rates happen to be equal. Then we add an infinitesimal weight to the string. The weight will slowly go down, and work will be done on the machine. Energy will be taken from the wheel and given to the vanes. If instead we take off a little bit of weight, then the imbalance is the other way. The weight is lifted, and heat is taken from the vane and put into the wheel. So we have the conditions of Carnot's reversible cycle, provided that the weight is just such that these two are equal. This condition is evidently that $(\epsilon + L\theta)/T_1 = \epsilon/T_2$. Let us say that the machine is slowly lifting the weight. Energy Q_1 is taken from the vanes and energy Q_2 is delivered to the wheel, and these energies are in the ratio $(\epsilon + L\theta)/\epsilon$. If we are lowering the weight, we also have $Q_1/Q_2 = (\epsilon + L\theta)/\epsilon$. Thus (Table 46-1) we have

$$Q_1/Q_2 = T_1/T_2.$$

Furthermore, the work we get out is to the energy taken from the vane as $L\theta$ is to $L\theta + \epsilon$, hence as $(T_1 - T_2)/T_1$. We see that our device cannot extract more work than this, operating reversibly. This is the result that we expected from Carnot's argument, and the main result of this lecture. However, we can use our device to understand a few other phenomena, even out of equilibrium, and therefore beyond the range of thermodynamics.

Let us now calculate *how fast* our one-way device would turn if everything were at the same temperature and we hung a weight on the drum. If we pull very, very hard, of course, there are all kinds of complications. The pawl slips over the ratchet, or the spring breaks, or something. But suppose we pull gently enough that everything works nicely. In those circumstances, the above analysis is right for the probability of the wheel going forward and backward, if we remember

that the two temperatures are equal. In each step an angle θ is obtained, so the angular velocity is θ times the probability of one of these jumps per second. It goes forward with probability $(1/\tau)e^{-(\epsilon+L\theta)/kT}$ and backward with probability $(1/\tau)e^{-\epsilon/kT}$, so that for the angular velocity we have

$$\begin{aligned}\omega &= (\theta/\tau)e^{-(\epsilon+L\theta)/kT} - e^{-\epsilon/kT} \\ &= (\theta/\tau)e^{-\epsilon/kT}(e^{-L\theta/kT} - 1).\end{aligned}\quad (46.1)$$

If we plot ω against L , we get the curve shown in Fig. 46-2. We see that it makes a great difference whether L is positive or negative. If L increases in the positive range, which happens when we try to drive the wheel backward, the backward velocity approaches a constant. As L becomes negative, ω really "takes off" forward, since e to a tremendous power is very great!

The angular velocity that was obtained from different forces is thus very unsymmetrical. Going one way it is easy: we get a lot of angular velocity for a little force. Going the other way, we can put on a lot of force, and yet the wheel hardly goes around.

We find the same thing in an *electrical rectifier*. Instead of the force, we have the electric field, and instead of the angular velocity, we have the electric current. In the case of a rectifier, the voltage is not proportional to resistance, and the situation is unsymmetrical. The same analysis that we made for the mechanical rectifier will also work for an electrical rectifier. In fact, the kind of formula we obtained above is typical of the current-carrying capacities of rectifiers as a function of their voltages.

Now let us take all the weights away, and look at the original machine. If T_2 were less than T_1 , the ratchet would go forward, as anybody would believe. But what is hard to believe, at first sight, is the opposite. If T_2 is greater than T_1 , the ratchet goes around the opposite way! A dynamic ratchet with lots of heat in it runs itself backwards, because the ratchet pawl is bouncing. If the pawl, for a moment, is on the incline somewhere, it pushes the inclined plane sideways. But it is *always* pushing on an inclined plane, because if it happens to lift up high enough to get past the point of a tooth, then the inclined plane slides by, and it comes down again on an inclined plane. So a hot ratchet and pawl is ideally built to go around in a direction exactly opposite to that for which it was originally designed!

In spite of all our cleverness of lopsided design, if the two temperatures are exactly equal there is no more propensity to turn one way than the other. The moment we look at it, it may be turning one way or the other, but in the long run, it gets nowhere. The fact that it gets nowhere is really the fundamental deep principle on which all of thermodynamics is based.

46-3 Reversibility in mechanics

What deeper mechanical principle tells us that, in the long run, if the temperature is kept the same everywhere, our gadget will turn neither to the right nor to the left? We evidently have a fundamental proposition that there is no way to design a machine which, left to itself, will be more likely to be turning one way than the other after a long enough time. We must try to see how this follows from the laws of mechanics.

The laws of mechanics go something like this: the mass times the acceleration is the force, and the force on each particle is some complicated function of the positions of all the other particles. There are other situations in which forces depend on velocity, such as in magnetism, but let us not consider that now. We take a simpler case, such as gravity, where forces depend only on position. Now suppose that we have solved our set of equations and we have a certain motion $x(t)$ for each particle. In a complicated enough system, the solutions are very complicated, and what happens with time turns out to be very surprising. If we write down any arrangement we please for the particles, we will see this arrangement actually occur if we wait long enough! If we follow our solution for a long

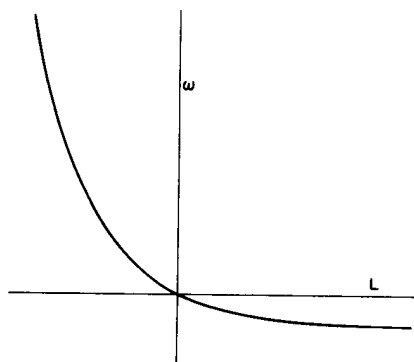


Fig. 46-2. Angular velocity of the ratchet as a function of torque.

enough time, it tries everything that it can do, so to speak. This is not absolutely necessary in the simplest devices, but when systems get complicated enough, with enough atoms, it happens. Now there is something else the solution can do. If we solve the equations of motion, we may get certain functions such as $t + t^2 + t^3$. We claim that another solution would be $-t + t^2 - t^3$. In other words, if we substitute $-t$ everywhere for t throughout the entire solution, we will once again get a solution of the same equation. This follows from the fact that if we substitute $-t$ for t in the original differential equation, nothing is changed, since only second derivatives with respect to t appear. This means that if we have a certain motion, then the exact opposite motion is also possible. In the complete confusion which comes if we wait long enough, it finds itself going one way sometimes, and it finds itself going the other way sometimes. There is nothing more beautiful about one of the motions than about the other. So it is impossible to design a machine which, in the long run, is more likely to be going one way than the other, if the machine is sufficiently complicated.

One might think up an example for which this is obviously untrue. If we take a wheel, for instance, and spin it in empty space, it will go the same way forever. So there are some conditions, like the conservation of angular momentum, which violate the above argument. This just requires that the argument be made with a little more care. Perhaps the walls take up the angular momentum, or something like that, so that we have no special conservation laws. Then, if the system is complicated enough, the argument is true. It is based on the fact that the laws of mechanics are reversible.

For historical interest, we would like to remark on a device invented by Maxwell, who first worked out the dynamical theory of gases. He supposed the following situation: We have two boxes of gas at the same temperature, with a little hole between them. At the hole sits a little demon (who may be a machine of course!). There is a door on the hole, which can be opened or closed by the demon. He watches the molecules coming from the left. Whenever he sees a fast molecule, he opens the door. When he sees a slow one, he leaves it closed. If we want him to be an extra special demon, he can have eyes at the back of his head, and do the opposite to the molecules from the other side. He lets the slow ones through to the left, and the fast through to the right. Pretty soon the left side will get cold and the right side hot. Then, are the ideas of thermodynamics violated because we could have such a demon?

It turns out, if we build a finite-sized demon, that the demon himself gets so warm that he cannot see very well after a while. The simplest possible demon, as an example, would be a trap door held over the hole by a spring. A fast molecule comes through, because it is able to lift the trap door. The slow molecule cannot get through, and bounces back. But this thing is nothing but our ratchet and pawl in another form, and ultimately the mechanism will heat up. If we assume that the specific heat of the demon is not infinite, it must heat up. It has but a finite number of internal gears and wheels, so it cannot get rid of the extra heat that it gets from observing the molecules. Soon it is shaking from Brownian motion so much that it cannot tell whether it is coming or going, much less whether the molecules are coming or going, so it does not work.

46-4 Irreversibility

Are all the laws of physics reversible? Evidently not! Just try to unscramble an egg! Run a moving picture backwards, and it takes only a few minutes for everybody to start laughing. The most natural characteristic of all phenomena is their obvious irreversibility.

Where does irreversibility come from? It does not come from Newton's laws. If we claim that the behavior of everything is ultimately to be understood in terms of the laws of physics, and if it also turns out that all the equations have the fantastic property that if we put $t = -t$ we have another solution, then every phenomenon is reversible. How then does it come about in nature on a large scale that things are not reversible? Obviously there must be some law, some

obscure but fundamental equation, perhaps in electricity, maybe in neutrino physics, in which it *does* matter which way time goes.

Let us discuss that question now. We already know one of those laws, which says that the entropy is always increasing. If we have a hot thing and a cold thing, the heat goes from hot to cold. So the law of entropy is one such law. But we expect to understand the law of entropy from the point of view of mechanics. In fact, we have just been successful in obtaining all the consequences of the argument that heat cannot flow backwards by itself from just mechanical arguments, and we thereby obtained an understanding of the Second Law. Apparently we can get irreversibility from reversible equations. But *was* it *only* a mechanical argument that we used? Let us look into it more closely.

Since our question has to do with the entropy, our problem is to try to find a microscopic description of entropy. If we say we have a certain amount of energy in something, like a gas, then we can get a microscopic picture of it, and say that every atom has a certain energy. All these energies added together give us the total energy. Similarly, maybe every atom has a certain entropy. If we add everything up, we would have the total entropy. It does not work so well, but let us see what happens.

As an example, we calculate the entropy difference between a gas at a certain temperature at one volume, and a gas at the same temperature at another volume. We remember, from Chapter 44, that we had, for the change in entropy,

$$\Delta S = \int \frac{dQ}{T}.$$

In the present case, the energy of the gas is the same before and after expansion, since the temperature does not change. So we have to add enough heat to equal the work done by the gas or, for each little change in volume,

$$dQ = P dV.$$

Putting this in for dQ , we get

$$\begin{aligned} \Delta S &= \int_{V_1}^{V_2} P \frac{dV}{T} = \int_{V_1}^{V_2} \frac{NkT}{V} \frac{dV}{T} \\ &= Nk \ln \frac{V_2}{V_1}, \end{aligned}$$

as we obtained in Chapter 44. For instance, if we expand the volume by a factor of 2, the entropy change is $Nk \ln 2$.

Let us now consider another interesting example. Suppose we have a box with a barrier in the middle. On one side is neon ("black" molecules), and on the other, argon ("white" molecules). Now we take out the barrier, and let them mix. How much has the entropy changed? It is possible to imagine that instead of the barrier we have a piston, with holes in it that let the whites through but not the blacks, and another kind of piston which is the other way around. If we move one piston to each end, we see that, for each gas, the problem is like the one we just solved. So we get an entropy change of $Nk \ln 2$, which means that the entropy has increased by $k \ln 2$ per molecule. The 2 has to do with the extra room that the molecule has, which is rather peculiar. It is not a property of the molecule itself, but of *how much room* the molecule has to run around in. This is a strange situation, where entropy increases but where everything has the same temperature and the same energy! The only thing that is changed is that the molecules are distributed differently.

We well know that if we just pull the barrier out, everything will get mixed up after a long time, due to the collisions, the jiggling, the banging, and so on. Every once in a while a white molecule goes toward a black, and a black one goes toward a white, and maybe they pass. Gradually the whites worm their way, by accident, across into the space of blacks, and the blacks worm their way, by accident, into the space of whites. If we wait long enough we get a mixture.

Clearly, this is an irreversible process in the real world, and ought to involve an increase in the entropy.

Here we have a simple example of an irreversible process which is completely composed of reversible events. Every time there is a collision between any two molecules, they go off in certain directions. If we took a moving picture of a collision in reverse, there would be nothing wrong with the picture. In fact, one kind of collision is just as likely as another. So the mixing is completely reversible, and yet it is irreversible. Everyone knows that if we started with white and with black, separated, we would get a mixture within a few minutes. If we sat and looked at it for several more minutes, it would not separate again but would stay mixed. So we have an irreversibility which is based on reversible situations. But we also see the *reason* now. We started with an arrangement which is, in some sense, *ordered*. Due to the chaos of the collisions, it becomes disordered. *It is the change from an ordered arrangement to a disordered arrangement which is the source of the irreversibility.*

It is true that if we took a motion picture of this, and showed it backwards, we would see it gradually become ordered. Someone would say, "That is against the laws of physics!" So we would run the film over again, and we would look at every collision. Every one would be perfect, and every one would be obeying the laws of physics. The reason, of course, is that every molecule's velocities are just right, so if the paths are all followed back, they get back to their original condition. But that is a very unlikely circumstance to have. If we start with the gas in no special arrangement, just whites and blacks, it will never get back.

46-5 Order and entropy

So we now have to talk about what we mean by disorder and what we mean by order. It is not a question of pleasant order or unpleasant disorder. What is different in our mixed and unmixed cases is the following. Suppose we divide the space into little volume elements. If we have white and black molecules, how many ways could we distribute them among the volume elements so that white is on one side, and black on the other? On the other hand, how many ways could we distribute them with no restriction on which goes where? Clearly, there are many more ways to arrange them in the latter case. We measure "disorder" by the number of ways that the insides can be arranged, so that from the outside it looks the same. *The logarithm of that number of ways is the entropy.* The number of ways in the separated case is less, so the entropy is less, or the "disorder" is less.

So with the above technical definition of disorder we can understand the proposition. First, the entropy measures the disorder. Second, the universe always goes from "order" to "disorder," so entropy always increases. Order is not order in the sense that we like the arrangement, but in the sense that the number of different ways we can hook it up, and still have it look the same from the outside, is relatively restricted. In the case where we reversed our motion picture of the gas mixing, there was not as much disorder as we thought. Every single atom had exactly the correct speed and direction to come out right! The entropy was not high after all, even though it appeared so.

What about the reversibility of the other physical laws? When we talked about the electric field which comes from an accelerating charge, it was said that we must take the retarded field. At a time t and at a distance r from the charge, we take the field due to the acceleration at a time $t - r/c$, not $t + r/c$. So it looks, at first, as if the law of electricity is not reversible. Very strangely, however, the laws we used come from a set of equations called Maxwell's equations, which are, in fact, reversible. Furthermore, it is possible to argue that if we were to use only the advanced field, the field due to the state of affairs at $t + r/c$, and do it absolutely consistently in a completely enclosed space, everything happens exactly the same way as if we use retarded fields! This apparent irreversibility in electricity, at least in an enclosure, is thus not an irreversibility at all. We have some feeling for that already, because we know that when we have an oscillating charge which generates fields which are bounced from the walls of an enclosure we ultimately

get to an equilibrium in which there is no one-sidedness. The retarded field approach is only a convenience in the method of solution.

So far as we know, all the fundamental laws of physics, like Newton's equations, are reversible. Then where does irreversibility come from? It comes from order going to disorder, but we do not understand this until we know the origin of the order. Why is it that the situations we find ourselves in every day are always out of equilibrium? One possible explanation is the following. Look again at our box of mixed white and black molecules. Now it is possible, if we wait long enough, by sheer, grossly improbable, but possible, accident, that the distribution of molecules gets to be mostly white on one side and mostly black on the other. After that, as times goes on and accidents continue, they get more mixed up again.

Thus one possible explanation of the high degree of order in the present-day world is that it is just a question of luck. Perhaps our universe happened to have had a fluctuation of some kind in the past, in which things got somewhat separated, and now they are running back together again. This kind of theory is not unsymmetrical, because we can ask what the separated gas looks like either a little in the future or a little in the past. In either case, we see a grey smear at the interface, because the molecules are mixing again. No matter which way we run time, the gas mixes. So this theory would say the irreversibility is just one of the accidents of life.

We would like to argue that this is not the case. Suppose we do not look at the whole box at once, but only at a piece of the box. Then, at a certain moment, suppose we discover a certain amount of order. In this little piece, white and black are separate. What should we deduce about the condition in places where we have not yet looked? If we really believe that the order arose from complete disorder by a fluctuation, we must surely take the most likely fluctuation which could produce it, and the most likely condition is *not* that the rest of it has also become disentangled! Therefore, from the hypothesis that the world is a fluctuation, all of the predictions are that if we look at a part of the world we have never seen before, we will find it mixed up, and not like the piece we just looked at. If our order were due to a fluctuation, we would not expect order anywhere but where we have just noticed it.

Now we assume the separation is because the past of the universe was really ordered. It is not due to a fluctuation, but the whole thing used to be white and black. This theory now predicts that there will be order in other places—the order is not due to a fluctuation, but due to a much higher ordering at the beginning of time. Then we would expect to find order in places where we have not yet looked.

The astronomers, for example, have only looked at some of the stars. Every day they turn their telescopes to other stars, and the new stars are doing the same thing as the other stars. We therefore conclude that the universe is *not* a fluctuation, and that the order is a memory of conditions when things started. This is not to say that we understand the logic of it. For some reason, the universe at one time had a very low entropy for its energy content, and since then the entropy has increased. So that is the way toward the future. That is the origin of all irreversibility, that is what makes the processes of growth and decay, that makes us remember the past and not the future, remember the things which are closer to that moment in the history of the universe when the order was higher than now, and why we are not able to remember things where the disorder is higher than now, which we call the future. So, as we commented in an earlier chapter, the entire universe is in a glass of wine, if we look at it closely enough. In this case the glass of wine is complex, because there is water and glass and light and everything else.

Another delight of our subject of physics is that even simple and idealized things, like the ratchet and pawl, work only because they are part of the universe. The ratchet and pawl works in only one direction because it has some ultimate contact with the rest of the universe. If the ratchet and pawl were in a box and isolated for some sufficient time, the wheel would be no more likely to go one way than the other. But because we pull up the shades and let the light out, because

we cool off on the earth and get heat from the sun, the ratchets and pawls that we make can turn one way. This one-wayness is interrelated with the fact that the ratchet is part of the universe. It is part of the universe not only in the sense that it obeys the physical laws of the universe, but its one-way behavior is tied to the one-way behavior of the entire universe. It cannot be completely understood until the mystery of the beginnings of the history of the universe are reduced still further from speculation to scientific understanding.

Sound. The wave equation

47-1 Waves

In this chapter we shall discuss the phenomenon of *waves*. This is a phenomenon which appears in many contexts throughout physics, and therefore our attention should be concentrated on it not only because of the particular example considered here, which is sound, but also because of the much wider application of the ideas in all branches of physics.

It was pointed out when we studied the harmonic oscillator that there are not only mechanical examples of oscillating systems but electrical ones as well. Waves are related to oscillating systems, except that wave oscillations appear not only as time-oscillations at one place, but propagate in space as well.

We have really already studied waves. When we studied light, in learning about the properties of waves in that subject, we paid particular attention to the interference in space of waves from several sources at different locations and all at the same frequency. There are two important wave phenomena that we have not yet discussed which occur in light, i.e., electromagnetic waves, as well as in any other form of waves. The first of these is the phenomenon of *interference in time* rather than interference in space. If we have two sources of sound which have slightly different frequencies and if we listen to both at the same time, then sometimes the waves come with the crests together and sometimes with the crest and trough together (see Fig. 47-1). The rising and falling of the sound that results is the phenomenon of *beats* or, in other words, of interference in time. The second phenomenon involves the wave patterns which result when the waves are confined within a given volume and reflect back and forth from walls.

These effects could have been discussed, of course, for the case of electromagnetic waves. The reason for not having done this is that by using one example we would not generate the feeling that we are actually learning about many different subjects at the same time. In order to emphasize the general applicability of waves beyond electrodynamics, we consider here a different example, in particular sound waves.

Other examples of waves are water waves consisting of long swells that we see coming in to the shore, or the smaller water waves consisting of surface tension ripples. As another example, there are two kinds of elastic waves in solids; a compressional (or longitudinal) wave in which the particles of the solid oscillate back and forth along the direction of propagation of the wave (sound waves in a gas are of this kind), and a transverse wave in which the particles of the solid oscillate in a direction perpendicular to the direction of propagation. Earthquake waves contain elastic waves of both kinds, generated by a motion at some place in the earth's crust.

Still another example of waves is found in modern physics. These are waves which give the probability amplitude of finding a particle at a given place—the “matter waves” which we have already discussed. Their frequency is proportional to the energy and their wave number is proportional to the momentum. They are the waves of quantum mechanics.

In this chapter we shall consider only waves for which the velocity is independent of the wavelength. This is, for example, the case for light in a vacuum. The speed of light is then the same for radiowaves, blue light, green light, or for any other wavelength. Because of this behavior, when we began to describe the wave phenomenon we did not notice at first that we had wave propagation. Instead, we said that if a charge is moved at one place, the electric field at a distance x was

47-1 Waves

47-2 The propagation of sound

47-3 The wave equation

47-4 Solutions of the wave equation

47-5 The speed of sound

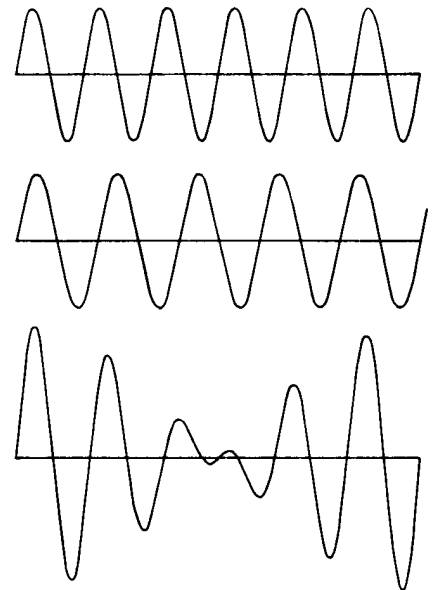


Fig. 47-1. Interference in time of two sound sources with slightly different frequencies, resulting in beats.

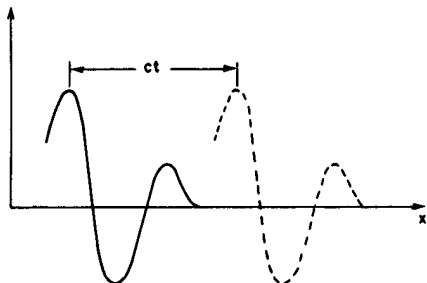


Fig. 47-2. The solid curve shows what the electric field might be like at some instant of time and the dashed curve shows what the electric field is at a time t later.

proportional to the acceleration, not at the time t , but at the earlier time $t - x/c$. Therefore if we were to picture the electric field in space at some instant of time, as in Fig. 47-2, the electric field at a time t later would have moved the distance ct , as indicated in the figure. Mathematically, we can say that in the one-dimensional example we are taking, the electric field is a function of $x - ct$. We see that at $t = 0$, it is some function of x . If we consider a later time, we need only to increase x somewhat to get the same value of the electric field. For example, if the maximum field occurred at $x = 3$ at time zero, then to find the new position of the maximum field at time t we need

$$x - ct = 3 \quad \text{or} \quad x = 3 + ct.$$

We see that this kind of function represents the propagation of a wave.

Such a function, $f(x - ct)$, then represents a wave. We may summarize this description of a wave by saying simply that

$$f(x - ct) = f(x + \Delta x - c(t + \Delta t)),$$

when $\Delta x = c \Delta t$. There is, of course, another possibility, i.e., that instead of a source to the left as indicated in Fig. 47-2, we have a source on the right, so that the wave propagates toward negative x . Then the wave would be described by $g(x + ct)$.

There is the additional possibility that more than one wave exists in space at the same time, and so the electric field is the sum of the two fields, each one propagating independently. This behavior of electric fields may be described by saying that if $f_1(x - ct)$ is a wave, and if $f_2(x - ct)$ is another wave, then their sum is also a wave. This is called the principle of superposition. The same principle is valid in sound.

We are familiar with the fact that if a sound is produced, we hear with complete fidelity the same sequence of sounds as was generated. If we had high frequencies travelling faster than low frequencies, a short, sharp noise would be heard as a succession of musical sounds. Similarly, if red light travelled faster than blue light, a flash of white light would be seen first as red, then as white, and finally as blue. We are familiar with the fact that this is not the case. Both sound and light travel with a speed in air which is very nearly independent of frequency. Examples of wave propagation for which this independence is not true will be considered in Chapter 48.

In the case of light (electromagnetic waves) we gave a rule which determined the electric field at a point as a result of the acceleration of a charge. One might expect now that what we should do is give a rule whereby some quality of the air, say the pressure, is determined at a given distance from a source in terms of the source motion, delayed by the travel time of the sound. In the case of light this procedure was acceptable because all that we knew was that a charge at one place exerts a force on another charge at another place. The details of propagation from the one place to the other were not absolutely essential. In the case of sound, however, we know that it propagates through the air between the source and the hearer, and it is certainly a natural question to ask what, at any given moment, the pressure of the air is. We would like, in addition, to know exactly how the air moves. In the case of electricity we could accept a rule, since we could say that we do not yet know the laws of electricity, but we cannot make the same remark with regard to sound. We would not be satisfied with a rule stating how the sound pressure moves through the air, because the process ought to be understandable as a consequence of the laws of mechanics. In short, sound is a branch of mechanics, and so it is to be understood in terms of Newton's laws. The propagation of sound from one place to another is merely a consequence of mechanics and the properties of gases, if it propagates in a gas, or of the properties of liquids or solids, if it propagates through such mediums. Later we shall derive the properties of light and its wave propagation in a similar way from the laws of electrodynamics.

47-2 The propagation of sound

We shall give a derivation of the properties of the propagation of sound *between* the source and the receiver as a consequence of Newton's laws, and we shall not consider the interaction with the source and the receiver. Ordinarily we emphasize a result rather than a particular derivation of it. In this chapter we take the opposite view. The point here, in a certain sense, is the derivation itself. This problem of explaining new phenomena in terms of old ones, when we know the laws of the old ones, is perhaps the greatest art of mathematical physics. The mathematical physicist has two problems: one is to find solutions, given the equations, and the other is to find the equations which describe a new phenomenon. The derivation here is an example of the second kind of problem.

We shall take the simplest example here—the propagation of sound in one dimension. To carry out such a derivation it is necessary first to have some kind of understanding of what is going on. Fundamentally what is involved is that if an object is moved at one place in the air, we observe that there is a disturbance which travels through the air. If we ask what kind of disturbance, we would say that we would expect that the motion of the object produces a change of pressure. Of course, if the object is moved gently, the air merely flows around it, but what we are concerned with is a rapid motion, so that there is not sufficient time for such a flow. Then, with the motion, the air is compressed and a change of pressure is produced which pushes on additional air. This air is in turn compressed, which leads again to an extra pressure, and a wave is propagated.

We now want to formulate such a process. We have to decide what variables we need. In our particular problem we would need to know how much the air has moved, so that the air *displacement* in the sound wave is certainly one relevant variable. In addition we would like to describe how the air *density* changes as it is displaced. The air *pressure* also changes, so this is another variable of interest. Then, of course, the air has a *velocity*, so that we shall have to describe the velocity of the air particles. The air particles also have *accelerations*—but as we list these many variables we soon realize that the velocity and acceleration would be known if we knew how the air *displacement* varies with time.

As we said, we shall consider the wave in one dimension. We can do this if we are sufficiently far from the source that what we call the *wavefronts* are very nearly planes. We thus make our argument simpler by taking the least complicated example. We shall then be able to say that the displacement, χ , depends only on x and t , and not on y and z . Therefore the description of the air is given by $\chi(x, t)$.

Is this description complete? It would appear to be far from complete, for we know none of the details of how the air molecules are moving. They are moving in all directions, and this state of affairs is certainly not described by means of this function $\chi(x, t)$. From the point of view of kinetic theory, if we have a higher density of molecules at one place and a lower density adjacent to that place, the molecules would move away from the region of higher density to the one of lower density, so as to equalize this difference. Apparently we would not get an oscillation and there would be no sound. What is necessary to get the sound wave is this situation: as the molecules rush out of the region of higher density and higher pressure, they give momentum to the molecules in the adjacent region of lower density. For sound to be generated, the regions over which the density and pressure change must be much larger than the distance the molecules travel before colliding with other molecules. This distance is the mean free path, and the distance between pressure crests and troughs must be much larger than this. Otherwise the molecules would move freely from the crest to the trough and immediately smear out the wave.

It is clear that we are going to describe the gas behavior on a scale large compared with the mean free path, and so the properties of the gas will not be described in terms of the individual molecules. The displacement, for example, will be the displacement of the center of mass of a small element of the gas, and the pressure or density will be the pressure or density in this region. We shall call the pressure P and the density ρ , and they will be functions of x and t . We must keep in mind that this description is an approximation which is valid only when these gas properties do not vary too rapidly with distance.

47-3 The wave equation

The physics of the phenomenon of sound waves thus involves three features:

- I. The gas moves and changes the density.
- II. The change in density corresponds to a change in pressure.
- III. Pressure inequalities generate gas motion.

Let us consider II first. For a gas, a liquid, or a solid, the pressure is some function of the density. Before the sound wave arrives, we have equilibrium, with a pressure P_0 and a corresponding density ρ_0 . A pressure P in the medium is connected to the density by some characteristic relation $P = f(\rho)$ and, in particular, the equilibrium pressure P_0 is given by $P_0 = f(\rho_0)$. The changes of pressure in sound from the equilibrium value are extremely small. A convenient unit for measuring pressure is the *bar*, where $1 \text{ bar} = 10^5 \text{ n/m}^2$. The pressure of 1 standard atmosphere is very nearly 1 bar: $1 \text{ atm} = 1.0133 \text{ bars}$. In sound we use a logarithmic scale of intensities since the sensitivity of the ear is roughly logarithmic. This scale is the decibel scale, in which the acoustic pressure level for the pressure amplitude P is defined as

$$I (\text{acoustic pressure level}) = 20 \log_{10}(P/P_{\text{ref}}) \text{ in db}, \quad (47.1)$$

where the reference pressure $P_{\text{ref}} = 2 \times 10^{-10} \text{ bar}$. A pressure amplitude of $P = 10^3 P_{\text{ref}} = 2 \times 10^{-7} \text{ bar}^*$ corresponds to a moderately intense sound of 60 decibels. We see that the pressure changes in sound are extremely small compared with the equilibrium, or mean, pressure of 1 atm. The displacements and the density changes are correspondingly extremely small. In explosions we do not have such small changes; the excess pressures produced can be greater than 1 atm. These large pressure changes lead to new effects which we shall consider later. In sound we do not often consider acoustic intensity levels over 100 db; 120 db is a level which is painful to the ear. Therefore, for sound, if we write

$$P = P_0 + P_e, \quad \rho = \rho_0 + \rho_e, \quad (47.2)$$

we shall always have the pressure change P_e very small compared with P_0 and the density change ρ_e very small compared with ρ_0 . Then

$$P_0 + P_e = f(\rho_0 + \rho_e) = f(\rho_0) + \rho_e f'(\rho_0), \quad (47.3)$$

where $P_0 = f(\rho_0)$ and $f'(\rho_0)$ stands for the derivative of $f(\rho)$ evaluated at $\rho = \rho_0$. We can take the second step in this equality only because ρ_e is very small. We find in this way that the excess pressure P_e is proportional to the excess density ρ_e , and we may call the proportionality factor κ :

$$P_e = \kappa \rho_e, \quad \text{where } \kappa = f'(\rho_0) = (dP/d\rho)_0. \quad (\text{II}) \quad (47.4)$$

The relation we needed for II is this very simple one.

Let us now consider I. We shall suppose that the position of a portion of air undisturbed by the sound wave is x and the displacement at the time t due to the sound is $\chi(x, t)$, so that its new position is $x + \chi(x, t)$, as in Fig. 47-3. Now the undisturbed position of a nearby portion of air is $x + \Delta x$, and its new position is $x + \Delta x + \chi(x + \Delta x, t)$. We can now find the density changes in the following way. Since we are limiting ourselves to plane waves, we can take a unit area perpendicular to the x -direction, which is the direction of propagation of the sound wave. The amount of air, per unit area, in Δx is then $\rho_0 \Delta x$, where ρ_0 is the undisturbed, or equilibrium, air density. This air, when displaced by the sound wave, now lies between $x + \chi(x, t)$ and $x + \Delta x + \chi(x + \Delta x, t)$, so that we have the same matter in this interval that was in Δx when undisturbed. If ρ is the new density, then

$$\rho_0 \Delta x = \rho [x + \Delta x + \chi(x + \Delta x, t) - x - \chi(x, t)]. \quad (47.5)$$

* With this choice of P_{ref} , the P is not the peak pressure in the sound wave but the "root-mean-square" pressure, which is $1/(2)^{1/2}$ times the peak pressure.

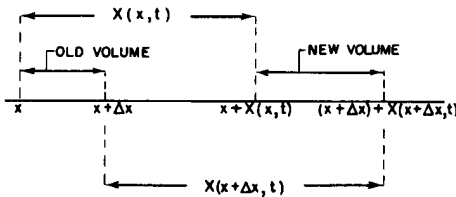


Fig. 47-3. The displacement of the air at x is $\chi(x, t)$, and at $x + \Delta x$ it is $\chi(x + \Delta x, t)$. The original volume of the air for a unit area of the plane wave is Δx ; the new volume is $\Delta x + \chi(x + \Delta x, t) - \chi(x, t)$.

Since Δx is small, we can write $\chi(x + \Delta x, t) - \chi(x, t) = (\partial\chi/\partial x) \Delta x$. This derivative is a partial derivative, since χ depends on the time as well as on x . Our equation then is

$$\rho_0 \Delta x = \rho \left(\frac{\partial\chi}{\partial x} \Delta x + \Delta x \right) \quad (47.6)$$

or

$$\rho_0 = (\rho_0 + \rho_e) \frac{\partial\chi}{\partial x} + \rho_0 + \rho_e. \quad (47.7)$$

Now in sound waves all changes are small so that ρ_e is small, χ is small, and $\partial\chi/\partial x$ is also small. Therefore in the relation that we have just found,

$$\rho_e = -\rho_0 \frac{\partial\chi}{\partial x} - \rho_e \frac{\partial\chi}{\partial x}, \quad (47.8)$$

we can neglect $\rho_e \partial\chi/\partial x$ compared with $\rho_0 \partial\chi/\partial x$. Thus we get the relation we needed for I:

$$\rho_e = -\rho_0 \frac{\partial\chi}{\partial x}. \quad (I) \quad (47.9)$$

This equation is what we would expect physically. If the displacements vary with x , then there will be density changes. The sign is also right: if the displacement χ increases with x , so that the air is stretched out, the density must go down.

We now need the third equation, which is the equation of the motion produced by the pressure. If we know the relation between the force and the pressure, we can then get the equation of motion. If we take a thin slab of air of length Δx and of unit area perpendicular to x , then the mass of air in this slab is $\rho_0 \Delta x$ and it has the acceleration $\partial^2\chi/\partial t^2$, so the mass times the acceleration for this slab of matter is $\rho_0 \Delta x (\partial^2\chi/\partial t^2)$. (It makes no difference for small Δx whether the acceleration $\partial^2\chi/\partial t^2$ is evaluated at an edge of the slab or at some intermediate position.) If now we find the force on this matter for a unit area perpendicular to x , it will then be equal to $\rho_0 \Delta x (\partial^2\chi/\partial t^2)$. We have the force in the $+x$ -direction, at x , of amount $P(x, t)$ per unit area, and we have the force in the opposite direction, at $x + \Delta x$, of amount $P(x + \Delta x, t)$ per unit area (Fig. 47-4):

$$P(x, t) - P(x + \Delta x, t) = -\frac{\partial P}{\partial x} \Delta x = -\frac{\partial P_e}{\partial x} \Delta x, \quad (47.10)$$

since Δx is small and since the only part of P which changes is the excess pressure P_e . We now have III:

$$\rho_0 \frac{\partial^2\chi}{\partial t^2} = -\frac{\partial P_e}{\partial x}, \quad (III) \quad (47.11)$$

and so we have enough equations to interconnect things and reduce down to one variable, say to χ . We can eliminate P_e from III by using II, so that we get

$$\rho_0 \frac{\partial^2\chi}{\partial t^2} = -\kappa \frac{\partial\rho_e}{\partial x}, \quad (47.12)$$

and then we can use I to eliminate ρ_e . In this way we find that ρ_0 cancels out and that we are left with

$$\frac{\partial^2\chi}{\partial t^2} = \kappa \frac{\partial^2\chi}{\partial x^2}. \quad (47.13)$$

We shall call $c_s^2 = \kappa$, so that we can write

$$\frac{\partial^2\chi}{\partial x^2} = \frac{1}{c_s^2} \frac{\partial^2\chi}{\partial t^2}. \quad (47.14)$$

This is the wave equation which describes the behavior of sound in matter.

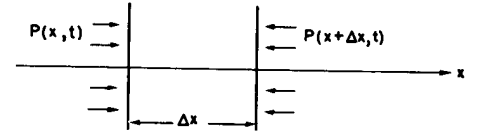


Fig. 47-4. The net force in the positive x -direction produced by the pressure acting on unit area perpendicular to x is $-(\partial P/\partial x) \Delta x$.

47-4 Solutions of the wave equation

We now can see whether this equation really does describe the essential properties of sound waves in matter. We want to deduce that a sound pulse, or disturbance, will move with a constant speed. We want to verify that two different pulses can move through each other—the principle of superposition. We also want to verify that sound can go either to the right or to the left. All these properties should be contained in this one equation.

We have remarked that any plane-wave disturbance which moves with a constant velocity v has the form $f(x - vt)$. Now we have to see whether $\chi(x, t) = f(x - vt)$ is a solution of the wave equation. When we calculate $\partial\chi/\partial x$, we get the derivative of the function, $\partial\chi/\partial x = f'(x - vt)$. Differentiating once more, we find

$$\frac{\partial^2 \chi}{\partial x^2} = f''(x - vt). \quad (47.15)$$

The differentiation of this same function with respect to t gives $-v$ times the derivative of the function, or $\partial\chi/\partial t = -vf'(x - vt)$, and the second time derivative is

$$\frac{\partial^2 \chi}{\partial t^2} = v^2 f''(x - vt). \quad (47.16)$$

It is evident that $f(x - vt)$ will satisfy the wave equation provided the wave velocity v is equal to c_s .

We find, therefore, from the *laws of mechanics* that any sound disturbance propagates with the velocity c_s , and in addition we find that

$$c_s = \kappa^{1/2} = (dP/d\rho)_0^{1/2},$$

and so we have related the wave velocity to a property of the medium.

If we consider a wave travelling in the opposite direction, so that $\chi(x, t) = g(x + vt)$, it is easy to see that such a disturbance also satisfies the wave equation. The only difference between such a wave and one travelling from left to right is in the sign of v , but whether we have $x + vt$ or $x - vt$ as the variable in the function does not affect the sign of $\partial^2 \chi / \partial t^2$, since it involves only v^2 . It follows that we have a solution for waves propagating in either direction with speed c_s .

An extremely interesting question is that of superposition. Suppose one solution of the wave equation has been found, say χ_1 . This means that the second derivative of χ_1 with respect to x is equal to $1/c_s^2$ times the second derivative of χ_1 with respect to t . Now any other solution χ_2 has this same property. If we superpose these two solutions, we have

$$\chi(x, t) = \chi_1(x, t) + \chi_2(x, t), \quad (47.17)$$

and we wish to verify that $\chi(x, t)$ is also a wave, i.e., that χ satisfies the wave equation. We can easily prove this result, since we have

$$\frac{\partial^2 \chi}{\partial x^2} = \frac{\partial^2 \chi_1}{\partial x^2} + \frac{\partial^2 \chi_2}{\partial x^2} \quad (47.18)$$

and, in addition,

$$\frac{\partial^2 \chi}{\partial t^2} = \frac{\partial^2 \chi_1}{\partial t^2} + \frac{\partial^2 \chi_2}{\partial t^2}. \quad (47.19)$$

It follows that $\partial^2 \chi / \partial x^2 = (1/c_s^2) \partial^2 \chi / \partial t^2$, so we have verified the principle of superposition. The proof of the principle of superposition follows from the fact that the wave equation is *linear* in χ .

We can now expect that a plane light wave propagating in the x -direction, polarized so that the electric field is in the y -direction, will satisfy the wave equation

$$\frac{\partial^2 E_y}{\partial x^2} = \frac{1}{c^2} \frac{\partial^2 E_y}{\partial t^2}, \quad (47.20)$$

where c is the speed of light. This wave equation is one of the consequences of Maxwell's equations. The equations of electrodynamics will lead to the wave equation for light just as the equations of mechanics lead to the wave equation for sound.

47-5 The speed of sound

Our deduction of the wave equation for sound has given us a *formula* which connects the wave speed with the rate of change of pressure with the density at the normal pressure:

$$c_s^2 = \left(\frac{dP}{d\rho} \right)_0. \quad (47.21)$$

In evaluating this rate of change, it is essential to know how the temperature varies. In a sound wave, we would expect that in the region of compression the temperature would be raised, and that in the region of rarefaction the temperature would be lowered. Newton was the first to calculate the rate of change of pressure with density, and he supposed that the temperature remained unchanged. He argued that the heat was conducted from one region to the other so rapidly that the temperature could not rise or fall. This argument gives the isothermal speed of sound, and it is wrong. The correct deduction was given later by Laplace, who put forward the opposite idea—that the pressure and temperature change adiabatically in a sound wave. The heat flow from the compressed region to the rarefied region is negligible so long as the wavelength is long compared with the mean free path. Under this condition the slight amount of heat flow in a sound wave does not affect the speed, although it gives a small absorption of the sound energy. We can expect correctly that this absorption increases as the wavelength approaches the mean free path, but these wavelengths are smaller by factors of about a million than the wavelengths of audible sound.

The actual variation of pressure with density in a sound wave is the one that allows no heat flow. This corresponds to the adiabatic variation, which we found to be $PV^\gamma = \text{const}$, where V was the volume. Since the density ρ varies inversely with V , the adiabatic connection between P and ρ is

$$P = \text{const } \rho^\gamma, \quad (47.22)$$

from which we get $dP/d\rho = \gamma P/\rho$. We then have for the speed of sound the relation

$$c_s^2 = \frac{\gamma P}{\rho}. \quad (47.23)$$

We can also write $c_s^2 = \gamma PV/\rho V$ and make use of the relation $PV = NkT$. Further, we see that ρV is the mass of gas, which can also be expressed as Nm , or as μ , where m is the mass of a molecule and μ is the molecular weight. In this way we find that

$$c_s^2 = \frac{\gamma kT}{m} = \frac{\gamma RT}{\mu}, \quad (47.24)$$

from which it is evident that the speed of sound depends only on the gas temperature and not on the pressure or the density. We also have observed that

$$kT = \frac{1}{3}m\langle v^2 \rangle, \quad (47.25)$$

where $\langle v^2 \rangle$ is the mean square of the speed of the molecules. It follows that $c_s^2 = (\gamma/3)\langle v^2 \rangle$, or

$$c_s = \left(\frac{\gamma}{3} \right)^{1/2} v_{\text{av}}. \quad (47.26)$$

This equation states that the speed of sound is some number which is roughly $1/(3)^{1/2}$ times some average speed, v_{av} , of the molecules (the square root of the

mean square velocity). In other words, the speed of sound is of the same order of magnitude as the speed of the molecules, and is actually somewhat less than this average speed.

Of course we could expect such a result, because a disturbance like a change in pressure is, after all, propagated by the motion of the molecules. However, such an argument does not tell us the precise propagation speed; it could have turned out that sound was carried primarily by the fastest molecules, or by the slowest molecules. It is reasonable and satisfying that the speed of sound is roughly $\frac{1}{2}$ of the average molecular speed v_{av} .

Beats

48-1 Adding two waves

Some time ago we discussed in considerable detail the properties of light waves and their interference—that is, the effects of the superposition of two waves from different sources. In all these analyses we assumed that the frequencies of the sources were all the same. In this chapter we shall discuss some of the phenomena which result from the interference of two sources which have *different* frequencies.

It is easy to guess what is going to happen. Proceeding in the same way as we have done previously, suppose we have two equal oscillating sources of the same frequency whose phases are so adjusted, say, that the signals arrive in phase at some point P . At that point, if it is light, the light is very strong; if it is sound, it is very loud; or if it is electrons, many of them arrive. On the other hand, if the arriving signals were 180° out of phase, we would get no signal at P , because the net amplitude there is then a minimum. Now suppose that someone twists the “phase knob” of one of the sources and changes the phase at P back and forth, say, first making it 0° and then 180° , and so on. Of course, we would then find variations in the net signal strength. Now we also see that if the phase of one source is slowly changing relative to that of the other in a gradual, uniform manner, starting at zero, going up to ten, twenty, thirty, forty degrees, and so on, then what we would measure at P would be a series of strong and weak “pulsations,” because when the phase shifts through 360° the amplitude returns to a maximum. Of course, to say that one source is shifting its phase relative to another at a uniform rate is the same as saying that the number of oscillations per second is slightly different for the two.

So we know the answer: if we have two sources at slightly different frequencies we should find, as a net result, an oscillation with a slowly pulsating intensity. That is all there really is to the subject!

It is very easy to formulate this result mathematically also. Suppose, for example, that we have two waves, and that we do not worry for the moment about all the spatial relations, but simply analyze what arrives at P . From one source, let us say, we would have $\cos \omega_1 t$, and from the other source, $\cos \omega_2 t$, where the two ω 's are not exactly the same. Of course the amplitudes may not be the same, either, but we can solve the general problem later; let us first take the case where the amplitudes are equal. Then the total amplitude at P is the sum of these two cosines. If we plot the amplitudes of the waves against the time, as in Fig. 48-1, we see that

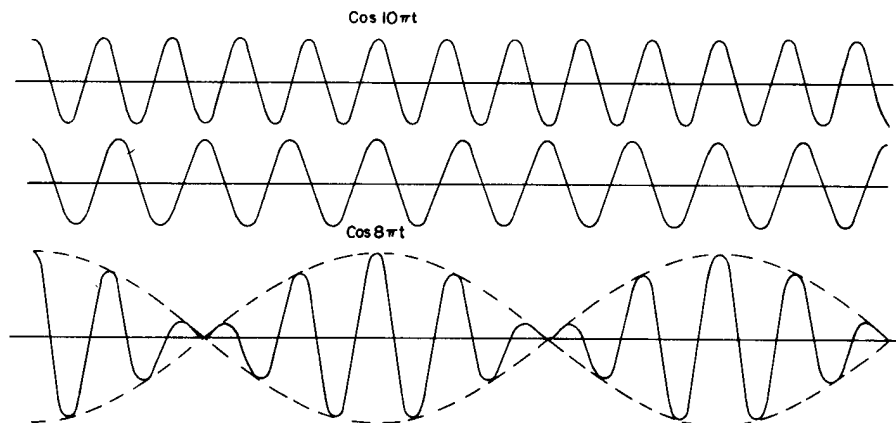


Fig. 48-1. The superposition of two cosine waves with frequencies in the ratio 8:10. The precise repetition of the pattern within each “beat” is not typical of the general case.

48-1 Adding two waves

48-2 Beat notes and modulation

48-3 Side bands

48-4 Localized wave trains

48-5 Probability amplitudes for particles

48-6 Waves in three dimensions

48-7 Normal modes

where the crests coincide we get a strong wave, and where a trough and crest coincide we get practically zero, and then when the crests coincide again we get a strong wave again.

Mathematically, we need only to add two cosines and rearrange the result somehow. There exist a number of useful relations among cosines which are not difficult to derive. Of course we know that

$$e^{i(a+b)} = e^{ia}e^{ib}, \quad (48.1)$$

and that e^{ia} has a real part, $\cos a$, and an imaginary part, $i \sin a$. If we take the real part of $e^{i(a+b)}$, we get $\cos(a+b)$. If we multiply out:

$$e^{ia}e^{ib} = (\cos a + i \sin a)(\cos b + i \sin b),$$

we get $\cos a \cos b - \sin a \sin b$, plus some imaginary parts. But we now need only the real part, so we have

$$\cos(a+b) = \cos a \cos b - \sin a \sin b. \quad (48.2)$$

Now if we change the sign of b , since the cosine does not change sign while the sine does, the same equation, for negative b , is

$$\cos(a-b) = \cos a \cos b + \sin a \sin b. \quad (48.3)$$

If we add these two equations together, we lose the sines and we learn that the product of two cosines is half the cosine of the sum, plus half the cosine of the difference:

$$\cos a \cos b = \frac{1}{2} \cos(a+b) + \frac{1}{2} \cos(a-b). \quad (48.4)$$

Now we can also reverse the formula and find a formula for $\cos \alpha + \cos \beta$ if we simply let $\alpha = a+b$ and $\beta = a-b$. That is, $a = \frac{1}{2}(\alpha + \beta)$ and $b = \frac{1}{2}(\alpha - \beta)$, so that

$$\cos \alpha + \cos \beta = 2 \cos \frac{1}{2}(\alpha + \beta) \cos \frac{1}{2}(\alpha - \beta). \quad (48.5)$$

Now we can analyze our problem. The sum of $\cos \omega_1 t$ and $\cos \omega_2 t$ is

$$\cos \omega_1 t + \cos \omega_2 t = 2 \cos \frac{1}{2}(\omega_1 + \omega_2)t \cos \frac{1}{2}(\omega_1 - \omega_2)t. \quad (48.6)$$

Now let us suppose that the two frequencies are nearly the same, so that $\frac{1}{2}(\omega_1 + \omega_2)$ is the average frequency, and is more or less the same as either. But $\omega_1 - \omega_2$ is *much smaller* than ω_1 or ω_2 because, as we suppose, ω_1 and ω_2 are nearly equal. That means that we can represent the solution by saying that there is a high-frequency cosine wave more or less like the ones we started with, but that its "size" is slowly changing—its "size" is pulsating with a frequency which appears to be $\frac{1}{2}(\omega_1 - \omega_2)$. But is this the frequency at which the beats are heard? Although (48.6) says that the amplitude goes as $\cos \frac{1}{2}(\omega_1 - \omega_2)t$, what it is really telling us is that the high-frequency oscillations are contained between two opposed cosine curves (shown dotted in Fig. 48-1). On this basis one could say that the amplitude varies at the frequency $\frac{1}{2}(\omega_1 - \omega_2)$, but if we are talking about the *intensity* of the wave we must think of it as having twice this frequency. That is, the modulation of the amplitude, in the sense of the strength of its intensity, is at frequency $\omega_1 - \omega_2$, although the formula tells us that we multiply by a cosine wave at half that frequency. The technical basis for the difference is that the high frequency-wave has a little different phase relationship in the second half-cycle.

Ignoring this small complication, we may conclude that if we add two waves of frequency ω_1 and ω_2 , we will get a net resulting wave of average frequency $\frac{1}{2}(\omega_1 + \omega_2)$ which oscillates in strength with a frequency $\omega_1 - \omega_2$.

If the two amplitudes are different, we can do it all over again by multiplying the cosines by different amplitudes A_1 and A_2 , and do a lot of mathematics, rearranging, and so on, using equations like (48.2)–(48.5). However, there are other, easier ways of doing the same analysis. For example, we know that it is much

easier to work with exponentials than with sines and cosines and that we can represent $A_1 \cos \omega_1 t$ as the real part of $A_1 e^{i\omega_1 t}$. The other wave would similarly be the real part of $A_2 e^{i\omega_2 t}$. If we add the two, we get $A_1 e^{i\omega_1 t} + A_2 e^{i\omega_2 t}$. If we then factor out the average frequency, we have

$$A_1 e^{i\omega_1 t} + A_2 e^{i\omega_2 t} = e^{1/2 i(\omega_1 + \omega_2)t} [A_1 e^{1/2 i(\omega_1 - \omega_2)t} + A_2 e^{-1/2 i(\omega_1 - \omega_2)t}]. \quad (48.7)$$

Again we have the high-frequency wave with a modulation at the lower frequency.

48-2 Beat notes and modulation

If we are now asked for the intensity of the wave of Eq. (48.7), we can either take the absolute square of the left side, or of the right side. Let us take the left side. The intensity then is

$$I = A_1^2 + A_2^2 + 2A_1 A_2 \cos(\omega_1 - \omega_2)t. \quad (48.8)$$

We see that the intensity swells and falls at a frequency $\omega_1 - \omega_2$, varying between the limits $(A_1 + A_2)^2$ and $(A_1 - A_2)^2$. If $A_1 \neq A_2$, the minimum intensity is not zero.

One more way to represent this idea is by means of a drawing, like Fig. 48-2. We draw a vector of length A_1 , rotating at a frequency ω_1 , to represent one of the waves in the complex plane. We draw another vector of length A_2 , going around at a frequency ω_2 , to represent the second wave. If the two frequencies are exactly equal, their resultant is of fixed length as it keeps revolving, and we get a definite, fixed intensity from the two. But if the frequencies are slightly different, the two complex vectors go around at different speeds. Figure 48-3 shows what the situation looks like relative to the vector $A_1 e^{i\omega_1 t}$. We see that A_2 is turning slowly away from A_1 , and so the amplitude that we get by adding the two is first strong, and then, as it opens out, when it gets to the 180° relative position the resultant gets particularly weak, and so on. As the vectors go around, the amplitude of the sum vector gets bigger and smaller, and the intensity thus pulsates. It is a relatively simple idea, and there are many different ways of representing the same thing.

The effect is very easy to observe experimentally. In the case of acoustics, we may arrange two loudspeakers driven by two separate oscillators, one for each loudspeaker, so that they each make a tone. We thus receive one note from one source and a different note from the other source. If we make the frequencies exactly the same, the resulting effect will have a definite strength at a given space location. If we then de-tune them a little bit, we hear some variations in the intensity. The farther they are de-tuned, the more rapid are the variations of sound. The ear has some trouble following variations more rapid than ten or so per second.

We may also see the effect on an oscilloscope which simply displays the sum of the currents to the two speakers. If the frequency of pulsing is relatively low, we simply see a sinusoidal wave train whose amplitude pulsates, but as we make the pulsations more rapid we see the kind of wave shown in Fig. 48-1. As we go to greater frequency differences, the "bumps" move closer together. Also, if the amplitudes are not equal and we make one signal stronger than the other, then we get a wave whose amplitude does not ever become zero, just as we expect. Everything works the way it should, both acoustically and electrically.

The opposite phenomenon occurs too! In radio transmission using so-called *amplitude modulation* (AM), the sound is broadcast by the radio station as follows: the radio transmitter has an AC electric oscillation which is at a very high frequency, for example 800 kilocycles per second, in the broadcast band. If this *carrier signal* is turned on, the radio station emits a wave which is of uniform amplitude at 800,000 oscillations a second. The way the "information" is transmitted, the useless kind of information about what kind of car to buy, is that when somebody talks into a microphone the amplitude of the carrier signal is changed in step with the vibrations of sound entering the microphone.

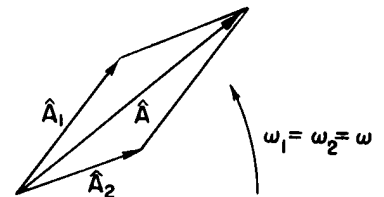


Fig. 48-2. The resultant of two complex vectors of equal frequency.

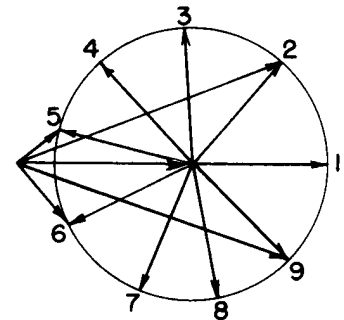


Fig. 48-3. The resultant of two complex vectors of unequal frequency, as seen in the rotating frame of reference of one vector. Nine successive positions of the slowly rotating vector are shown.

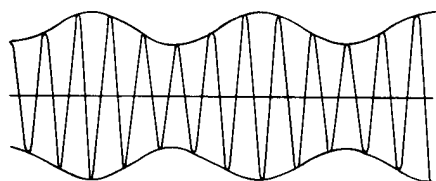


Fig. 48-4. A modulated carrier wave. In this schematic sketch, $\omega_c/\omega_m = 5$. In an actual radiowave, $\omega_c/\omega_m \sim 100$.

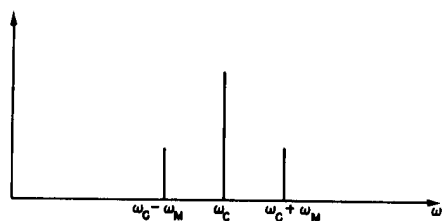


Fig. 48-5. The frequency spectrum of a carrier wave ω_c modulated by a single cosine wave ω_m .

If we take as the simplest mathematical case the situation where a soprano is singing a perfect note, with perfect sinusoidal oscillations of her vocal cords, then we get a signal whose strength is alternating as shown in Fig. 48-4. The audiofrequency alternation is then recovered in the receiver; we get rid of the carrier wave and just look at the envelope which represents the oscillations of the vocal cords, or the sound of the singer. The loudspeaker then makes corresponding vibrations at the same frequency in the air, and the listener is then essentially unable to tell the difference, so they say. Because of a number of distortions and other subtle effects, it is, in fact, possible to tell whether we are listening to a radio or to a real soprano; otherwise the idea is as indicated above.

48-3 Side bands

Mathematically, the modulated wave described above would be expressed as

$$S = (1 + b \cos \omega_m t) \cos \omega_c t, \quad (48.9)$$

where ω_c represents the frequency of the carrier and ω_m is the frequency of the audio tone. Again we use all those theorems about the cosines, or we can use $e^{i\theta}$; it makes no difference—it is easier with $e^{i\theta}$, but it is the same thing. We then get

$$S = \cos \omega_c t + \frac{1}{2}b \cos (\omega_c + \omega_m)t + \frac{1}{2}b \cos (\omega_c - \omega_m)t. \quad (48.10)$$

So, from another point of view, we can say that the output wave of the system consists of three waves added in superposition: first, the regular wave at the frequency ω_c , that is, at the carrier frequency, and then two new waves at two new frequencies. One is the carrier frequency plus the modulation frequency, and the other is the carrier frequency minus the modulation frequency. If, therefore, we make some kind of plot of the intensity being generated by the generator as a function of frequency, we would find a lot of intensity at the frequency of the carrier, naturally, but when a singer started to sing, we would suddenly also find intensity proportional to the strength of the singer, b^2 , at frequency $\omega_c + \omega_m$ and $\omega_c - \omega_m$, as shown in Fig. 48-5. These are called *side bands*; when there is a modulated signal from the transmitter, there are side bands. If there is more than one note at the same time, say ω_m and ω_m' , there are two instruments playing; or if there is any other complicated cosine wave, then, of course, we can see from the mathematics that we get some more waves that correspond to the frequencies $\omega_c \pm \omega_m$.

Therefore, when there is a complicated modulation that can be represented as the sum of many cosines,* we find that the actual transmitter is transmitting over a range of frequencies, namely the carrier frequency plus or minus the maximum frequency that the modulation signal contains.

Although at first we might believe that a radio transmitter transmits only at the nominal frequency of the carrier, since there are big, superstable crystal oscillators in there, and everything is adjusted to be at precisely 800 kilocycles, the moment someone *announces* that they are at 800 kilocycles, he modulates the 800 kilocycles, and so they are no longer precisely at 800 kilocycles! Suppose that the amplifiers are so built that they are able to transmit over a good range of the ear's sensitivity (the ear can hear up to 20,000 cycles per second, but usually radio transmitters and receivers do not work beyond 10,000, so we do not hear the highest parts), then, when the man speaks, his voice may contain frequencies ranging up, say, to 10,000 cycles, so the transmitter is transmitting frequencies

* A slight side remark: In what circumstances can a curve be represented as a sum of a lot of cosines? *Answer:* In all ordinary circumstances, except for certain cases the mathematicians can dream up. Of course, the curve must have only one value at a given point, and it must not be a crazy curve which jumps an infinite number of times in an infinitesimal distance, or something like that. But aside from such restrictions any reasonable curve (one that a singer is going to be able to make by shaking her vocal cords) can always be compounded by adding cosine waves together.

which may range from 790 to 810 kilocycles per second. Now if there were another station at 795 kc/sec, there would be a lot of confusion. Also, if we made our receiver so sensitive that it picked up only 800, and did not pick up the 10 kilocycles on either side, we would not hear what the man was saying, because the information would be on these other frequencies! Therefore it is absolutely essential to keep the stations a certain distance apart, so that their side bands do not overlap and, also, the receiver must not be so selective that it does not permit reception of the side bands as well as of the main nominal frequency. In the case of sound, this problem does not really cause much trouble. We can hear over a ± 20 kc/sec range, and we have usually from 500 to 1500 kc/sec in the broadcast band, so there is plenty of room for lots of stations.

The television problem is more difficult. As the electron beam goes across the face of the picture tube, there are various little spots of light and dark. That "light" and "dark" is the "signal." Now ordinarily the beam scans over the whole picture, 500 lines, approximately, in a thirtieth of a second. Let us consider that the resolution of the picture vertically and horizontally is more or less the same, so that there are the same number of spots per inch along a scan line. We want to be able to distinguish dark from light, dark from light, dark from light, over, say, 500 lines. In order to be able to do this with cosine waves, the shortest wavelength needed thus corresponds to a wavelength, from maximum to maximum, of one 250th of the screen size. So we have $250 \times 500 \times 30$ pieces of information per second. The highest frequency that we are going to carry, therefore, is close to 4 megacycles per second. Actually, to keep the television stations apart, we have to use a little bit more than this, about 6 mc/sec; part of it is used to carry the sound signal, and other information. So, television channels are 6 megacycles per second wide. It certainly would not be possible to transmit TV on an 800 kc/sec carrier, since we cannot modulate at a higher frequency than the carrier.

At any rate, the television band starts at 54 megacycles. The first transmission channel, which is channel 2 (!), has a frequency range from 54 to 60 mc/sec, which is 6 mc/sec wide. "But," one might say, "we have just proved that there were side bands on both sides, and therefore it should be twice that wide." It turns out that the radio engineers are rather clever. If we analyze the modulation signal using not just cosine terms, but cosine and sine terms, to allow for phase differences, we then see that there is a definite, invariant relationship between the side band on the high-frequency side and the side band on the low-frequency side. What we mean is that there is no new information on that other side band. So what is done is to suppress one side band, and the receiver is wired inside such that the information which is missing is reconstituted by looking at the single side band and the carrier. Single side-band transmission is a clever scheme for decreasing the band widths needed to transmit information.

48-4 Localized wave trains

The next subject we shall discuss is the interference of waves in both space and time. Suppose that we have two waves travelling in space. We know, of course, that we can represent a wave travelling in space by $e^{i(\omega t - kx)}$. This might be, for example, the displacement in a sound wave. This is a solution of the wave equation provided that $\omega^2 = k^2 c^2$, where c is the speed of propagation of the wave. In this case we can write it as $e^{ik(x - ct)}$, which is of the general form $f(x - ct)$. Therefore this must be a wave which is travelling at this velocity, ω/k , and that is c and everything is all right.

Now we want to add two such waves together. Suppose we have a wave that is travelling with one frequency, and another wave travelling with another frequency. We leave to the reader to consider the case where the amplitudes are different; it makes no real difference. Thus we want to add $e^{i(\omega_1 t - k_1 x)} + e^{i(\omega_2 t - k_2 x)}$. We can add these by the same kind of mathematics we used when we added signal waves. Of course, if c is the same for both, this is easy, since it is the same as what we did before:

$$e^{i\omega_1(t-x/c)} + e^{i\omega_2(t-x/c)} = e^{i\omega_1 t'} + e^{i\omega_2 t'}, \quad (48.11)$$

except that $t' = t - x/c$ is the variable instead of t . So we get the same kind of modulations, naturally, but we see, of course, that those modulations are moving along with the wave. In other words, if we added two waves, but these waves were not just oscillating, but also moving in space, then the resultant wave would move along also, at the same speed.

Now we would like to generalize this to the case of waves in which the relationship between the frequency and the wave number k is not so simple. Example: material having an index of refraction. We have already studied the theory of the index of refraction in Chapter 31, where we found that we could write $k = n\omega/c$, where n is the index of refraction. As an interesting example, for x-rays we found that the index n is

$$n = 1 - \frac{Nq_e^2}{2\epsilon_0 m \omega^2}. \quad (48.12)$$

We actually derived a more complicated formula in Chapter 31, but this one is as good as any, as an example.

Incidentally, we know that even when ω and k are not linearly proportional, the ratio ω/k is certainly the speed of propagation for the particular frequency and wave number. We call this ratio the *phase velocity*; it is the speed at which the phase, or the nodes of a single wave, would move along:

$$v_p = \frac{\omega}{k}. \quad (48.13)$$

This phase velocity, for the case of x-rays in glass, is greater than the speed of light in vacuum (since n in 48.12 is less than 1), and that is a bit bothersome, because we do not think we can send signals faster than the speed of light!

What we are going to discuss now is the interference of two waves in which ω and k have a definite formula relating them. The above formula for n says that k is given as a definite function of ω . To be specific, in this particular problem, the formula for k in terms of ω is

$$k = \frac{\omega}{c} - \frac{a}{\omega c}, \quad (48.14)$$

where $a = Nq_e^2/2\epsilon_0 m$, a constant. At any rate, for each frequency there is a definite wave number, and we want to add two such waves together.

Let us do it just as we did in Eq. (48.7):

$$e^{i(\omega_1 t - k_1 x)} + e^{i(\omega_2 t - k_2 x)} = e^{1/2 i[(\omega_1 + \omega_2)t - (k_1 + k_2)x]} \times \{e^{1/2 i[(\omega_1 - \omega_2)t - (k_1 - k_2)x]} + e^{-1/2 i[(\omega_1 - \omega_2)t - (k_1 - k_2)x]}\}. \quad (48.15)$$

So we have a modulated wave again, a wave which travels with the mean frequency and the mean wave number, but whose strength is varying with a form which depends on the difference frequency and the difference wave number.

Now let us take the case that the difference between the two waves is relatively small. Let us suppose that we are adding two waves whose frequencies are nearly equal; then $(\omega_1 + \omega_2)/2$ is practically the same as either one of the ω 's, and similarly for $(k_1 + k_2)/2$. Thus the speed of the wave, the fast oscillations, the nodes, is still essentially ω/k . But look, the speed of propagation of the modulation is not the same! How much do we have to change x to account for a certain amount of t ? The speed of this modulation wave is the ratio

$$v_M = \frac{\omega_1 - \omega_2}{k_1 - k_2}. \quad (48.16)$$

The speed of modulation is sometimes called the *group velocity*. If we take the case that the difference in frequency is relatively small, and the difference in wave number is then also relatively small, then this expression approaches, in the limit,

$$v_g = \frac{d\omega}{dk}. \quad (48.17)$$

In other words, for the slowest modulation, the slowest beats, there is a definite speed at which they travel which is not the same as the phase speed of the waves—what a mysterious thing!

The group velocity is the derivative of ω with respect to k , and the phase velocity is ω/k .

Let us see if we can understand why. Consider two waves, again of slightly different wavelength, as in Fig. 48-1. They are out of phase, in phase, out of phase, and so on. Now these waves represent, really, the waves in space travelling with slightly different frequencies also. Now because the phase velocity, the velocity of the nodes of these two waves, is not precisely the same, something new happens. Suppose we ride along with one of the waves and look at the other one; if they both went at the same speed, then the other wave would stay right where it was relative to us, as we ride along on this crest. We ride on that crest and right opposite us we see a crest; if the two velocities are equal the crests stay on top of each other. But it is *not so* that the two velocities are really equal. There is only a small difference in frequency and therefore only a small difference in velocity, but because of that difference in velocity, as we ride along the other wave moves slowly forward, say, or behind, relative to our wave. So as time goes on, what happens to the node? If we move one wave train just a shade forward, the node moves forward (or backward) a considerable distance. That is, the sum of these two waves has an envelope, and as the waves travel along, the envelope rides on them at a different speed. The *group velocity* is the speed at which modulated signals would be transmitted.

If we made a signal, i.e., some kind of change in the wave that one could recognize when he listened to it, a kind of modulation, then that modulation would travel at the group velocity, provided that the modulations were relatively slow. (When they are fast, it is much more difficult to analyze.)

Now we may show (at long last), that the speed of propagation of x-rays in a block of carbon is *not* greater than the speed of light, although the phase velocity *is* greater than the speed of light. In order to do that, we must find $d\omega/dk$, which we get by differentiating (48.14): $dk/d\omega = 1/c + a/\omega^2 c$. The group velocity, therefore, is the reciprocal of this, namely,

$$v_g = \frac{c}{1 + a/\omega^2}, \quad (48.18)$$

which is smaller than c ! So although the phases can travel faster than the speed of light, the modulation signals travel slower, and that is the resolution of the apparent paradox! Of course, if we have the simple case that $\omega = kc$, then $d\omega/dk$ is also c . So when all the phases have the same velocity, naturally the group has the same velocity.

48-5 Probability amplitudes for particles

Let us now consider one more example of the phase velocity which is extremely interesting. It has to do with quantum mechanics. We know that the amplitude to find a particle at a place can, in some circumstances, vary in space and time, let us say in one dimension, in this manner:

$$\psi = Ae^{i(\omega t - kx)}, \quad (48.19)$$

where ω is the frequency, which is related to the classical idea of the energy through $E = \hbar\omega$, and k is the wave number, which is related to the momentum through $p = \hbar k$. We would say the particle had a definite momentum p if the wave number were exactly k , that is, a perfect wave which goes on with the same amplitude everywhere. Equation (48.19) gives the amplitude, and if we take the absolute square, we get the relative probability for finding the particle as a function of position and time. This is a *constant*, which means that the probability is the same to find a particle anywhere. Now suppose, instead, that we have a situation where we know that the particle is more likely to be at one place than at another. We would

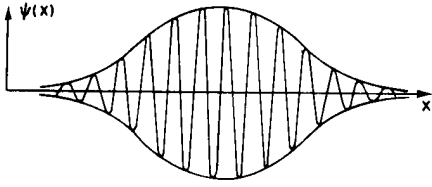


Fig. 48-6. A localized wave train.

represent such a situation by a wave which has a maximum and dies out on either side (Fig. 48-6). (It is not quite the same as a wave like (48.1) which has a series of maxima, but it is possible, by adding several waves of nearly the same ω and k together, to get rid of all but one maximum.)

Now in those circumstances, since the square of (48.19) represents the chance of finding a particle somewhere, we know that at a given instant the particle is most likely to be near the center of the "lump," where the amplitude of the wave is maximum. If now we wait a few moments, the waves will move, and after some time the "lump" will be somewhere else. If we knew that the particle originally was situated somewhere, classically, we would *expect* that it would later be elsewhere as a matter of fact, because it has a *speed*, after all, and a momentum. The quantum theory, then, will go into the correct classical theory for the relationship of momentum, energy, and velocity only if the group velocity, the velocity of the modulation, is equal to the velocity that we would obtain classically for a particle of the same momentum.

It is now necessary to demonstrate that this is, or is not, the case. According to the classical theory, the energy is related to the velocity through an equation like

$$E = \frac{mc^2}{\sqrt{1 - v^2/c^2}}. \quad (48.20)$$

Similarly, the momentum is

$$p = \frac{mv}{\sqrt{1 - v^2/c^2}}. \quad (48.21)$$

That is the classical theory, and as a consequence of the classical theory, by eliminating v , we can show that

$$E^2 - p^2c^2 = m^2c^4.$$

That is the four-dimensional grand result that we have talked and talked about, that $p_\mu p_\mu = m^2$; that is the relation between energy and momentum in the classical theory. Now that means, since these E 's and p 's are going to become ω 's and k 's, by substitution of $E = \hbar\omega$ and $p = \hbar k$, that for quantum mechanics it is necessary that

$$\frac{\hbar^2\omega^2}{c^2} - \hbar^2k^2 = m^2c^2. \quad (48.22)$$

This, then, is the relationship between the frequency and the wave number of a quantum-mechanical amplitude wave representing a particle of mass m . From this equation we can deduce that ω is

$$\omega = c\sqrt{k^2 + m^2c^2/\hbar^2}.$$

The phase velocity, ω/k , is here again faster than the speed of light!

Now let us look at the group velocity. The group velocity should be $d\omega/dk$, the speed at which the modulations move. We have to differentiate a square root, which is not very difficult. The derivative is

$$\frac{d\omega}{dk} = \frac{kc}{\sqrt{k^2 + m^2c^2/\hbar^2}}.$$

Now the square root is, after all, ω , so we could write this as $d\omega/dk = c^2k/\omega$. Further, k/ω is p/E , so

$$v_g = \frac{c^2p}{E}.$$

But from (48.20) and (48.21), $c^2p/E = v$, the velocity of the particle, according to classical mechanics. So we see that whereas the fundamental quantum-mechanical relationship $E = \hbar\omega$ and $p = \hbar k$, for the identification of ω and k with the classical E and p , only produces the equation $\omega^2 - k^2c^2 = m^2c^4/\hbar^2$, now we also understand the relationships (48.20) and (48.21) which connected E and p

to the velocity. Of course the group velocity must be the velocity of the particle if the interpretation is going to make any sense. If we think the particle is over here at one time, and then ten minutes later we think it is over there, as the quantum mechanics said, the distance traversed by the "lump," divided by the time interval, must be, classically, the velocity of the particle.

48-6 Waves in three dimensions

We shall now bring our discussion of waves to a close with a few general remarks about the wave equation. These remarks are intended to give some view of the future—not that we can understand everything exactly just now, but rather to see what things are going to look like when we study waves a little more. First of all, the wave equation for sound in one dimension was

$$\frac{\partial^2 \chi}{\partial x^2} = \frac{1}{c^2} \frac{\partial^2 \chi}{\partial t^2},$$

where c is the speed of whatever the wave is—in the case of sound, it is the sound speed; in the case of light, it is the speed of light. We showed that for a sound wave the displacements would propagate themselves at a certain speed. But the excess pressure also propagates at a certain speed, and so does the excess density. So we should expect that the pressure would satisfy the same equation, as indeed it does. We shall leave it to the reader to prove that it does. *Hint:* ρ_e is proportional to the rate of change of χ with respect to x . Therefore if we differentiate the wave equation with respect to x , we will immediately discover that $\partial \chi / \partial x$ satisfies the same equation. That is to say, ρ_e satisfies the same equation. But P_e is proportional to ρ_e , and therefore P_e does too. So the pressure, the displacements, everything, satisfy the same wave equation.

Usually one sees the wave equation for sound written in terms of pressure instead of in terms of displacement, because the pressure is a scalar and has no direction. But the displacement is a vector and has direction, and it is thus easier to analyze the pressure.

The next matter we discuss has to do with the wave equation in three dimensions. We know that the sound wave solution in one dimension is $e^{i(\omega t - kx)}$, with $\omega = kc_s$, but we also know that in three dimensions a wave would be represented by $e^{i(\omega t - k_x x - k_y y - k_z z)}$, where, in this case, $\omega^2 = k^2 c_s^2$, which is, of course, $(k_x^2 + k_y^2 + k_z^2) c_s^2$. Now what we want to do is to guess what the correct wave equation in three dimensions is. Naturally, for the case of sound this can be deduced by going through the same dynamic argument in three dimensions that we made in one dimension. But we shall not do that; instead we just write down what comes out: the equation for the pressure (or displacement, or anything) is

$$\frac{\partial^2 P_e}{\partial x^2} + \frac{\partial^2 P_e}{\partial y^2} + \frac{\partial^2 P_e}{\partial z^2} = \frac{1}{c_s^2} \frac{\partial^2 P_e}{\partial t^2}. \quad (48.23)$$

That this is true can be verified by substituting in $e^{i(\omega t - \mathbf{k} \cdot \mathbf{r})}$. Clearly, every time we differentiate with respect to x , we multiply by $-ik_x$. If we differentiate twice, it is equivalent to multiplying by $-k_x^2$, so the first term would become $-k_x^2 P_e$, for that wave. Similarly, the second term becomes $-k_y^2 P_e$, and the third term becomes $-k_z^2 P_e$. On the right, we get $-(\omega^2/c_s^2) P_e$. Then, if we take away the P_e 's and change the sign, we see that the relationship between k and ω is the one that we want.

Working backwards again, we cannot resist writing down the grand equation which corresponds to the dispersion equation (48.22) for quantum-mechanical waves. If ϕ represents the amplitude for finding a particle at position x, y, z , at the time t , then the great equation of quantum mechanics for free particles is this:

$$\frac{\partial^2 \phi}{\partial x^2} + \frac{\partial^2 \phi}{\partial y^2} + \frac{\partial^2 \phi}{\partial z^2} - \frac{1}{c^2} \frac{\partial^2 \phi}{\partial t^2} = \frac{m^2 c^2}{\hbar^2} \phi. \quad (48.24)$$

First of all, the relativity character of this expression is suggested by the appearance

of x , y , z and t in the nice combination relativity usually involves. Second, it is a wave equation which, if we try a plane wave, would produce as a consequence that $-k^2 + \omega^2/c^2 = m^2c^2/\hbar^2$, which is the right relationship for quantum mechanics. There is still another great thing contained in the wave equation: the fact that any superposition of waves is also a solution. So this equation contains all of the quantum mechanics and the relativity that we have been discussing so far, at least so long as it deals with a single particle in empty space with no external potentials or forces on it!

48-7 Normal modes

Now we turn to another example of the phenomenon of beats which is rather curious and a little different. Imagine two equal pendulums which have, between them, a rather weak spring connection. They are made as nearly as possible the same length. If we pull one aside and let go, it moves back and forth, and it pulls on the connecting spring as it moves back and forth, and so it really is a machine for generating a force which has the natural frequency of the other pendulum. Therefore, as a consequence of the theory of resonance, which we studied before, when we put a force on something at just the right frequency, it will drive it. So, sure enough, one pendulum moving back and forth drives the other. However, in this circumstance there is a new thing happening, because the total energy of the system is finite, so when one pendulum pours its energy into the other to drive it, it finds itself gradually losing energy, until, if the timing is just right along with the speed, it loses all its energy and is reduced to a stationary condition! Then, of course, it is the other pendulum ball that has all the energy and the first one which has none, and as time goes on we see that it works also in the opposite direction, and that the energy is passed back into the first ball; this is a very interesting and amusing phenomenon. We said, however, that this is related to the theory of beats, and we must now explain how we can analyze this motion from the point of view of the theory of beats.

We note that the motion of either of the two balls is an oscillation which has an amplitude which changes cyclically. Therefore the motion of one of the balls is presumably analyzable in a different way, in that it is the sum of two oscillations, present at the same time but having two slightly different frequencies. Therefore it ought to be possible to find two other motions in this system, and to claim that what we saw was a superposition of the two solutions, because this is of course a linear system. Indeed, it is easy to find two ways that we could start the motion, each one of which is a perfect, single-frequency motion—absolutely periodic. The motion that we started with before was not strictly periodic, since it did not last; soon one ball was passing energy to the other and so changing its amplitude; but there are ways of starting the motion so that nothing changes and, of course, as soon as we see it we understand why. For example, if we made both pendulums go together, then, since they are of the same length and the spring is not then doing anything, they will of course continue to swing like that for all time, assuming no friction and that everything is perfect. On the other hand, there is another possible motion which also has a definite frequency: that is, if we move the pendulums oppositely, pulling them aside exactly equal distances, then again they would be in absolutely periodic motion. We can appreciate that the spring just adds a little to the restoring force that the gravity supplies, that is all, and the system just keeps oscillating at a slightly higher frequency than in the first case. Why higher? Because the spring is pulling, in addition to the gravitation, and it makes the system a little “stiffer,” so that the frequency of this motion is just a shade higher than that of the other.

Thus this system has two ways in which it can oscillate with unchanging amplitude: it can either oscillate in a manner in which both pendulums go the same way and oscillate all the time at one frequency, or they could go in opposite directions at a slightly higher frequency.

Now the actual motion of the thing, because the system is linear, can be represented as a superposition of the two. (The subject of this chapter, remember,

is the effects of adding two motions with different frequencies.) So think what would happen if we combined these two solutions. If at $t = 0$ the two motions are started with equal amplitude and in the same phase, the sum of the two motions means that one ball, having been impressed one way by the first motion and the other way by the second motion, is at zero, while the other ball, having been displaced the same way in both motions, has a large amplitude. As time goes on, however, the two basic *motions* proceed independently, so the phase of one relative to the other is slowly shifting. That means, then, that after a sufficiently long time, when the time is enough that one motion could have gone “ $900\frac{1}{2}$ ” oscillations, while the other went only “900,” the relative phase would be just reversed with respect to what it was before. That is, the large-amplitude motion will have fallen to zero, and in the meantime, of course, the initially motionless ball will have attained full strength!

So we see that we could analyze this complicated motion either by the idea that there is a resonance and that one passes energy to the other, or else by the superposition of two constant-amplitude motions at two different frequencies.

Modes

49-1 The reflection of waves

This chapter will consider some of the remarkable phenomena which are a result of confining waves in some finite region. We will be led first to discover a few particular facts about vibrating strings, for example, and then the generalization of these facts will give us a principle which is probably the most far-reaching principle of mathematical physics.

Our first example of confining waves will be to confine a wave at one boundary. Let us take the simple example of a one-dimensional wave on a string. One could equally well consider sound in one dimension against a wall, or other situations of a similar nature, but the example of a string will be sufficient for our present purposes. Suppose that the string is held at one end, for example by fastening it to an "infinitely solid" wall. This can be expressed mathematically by saying that the displacement y of the string at the position $x = 0$ must be zero, because the end does not move. Now if it were not for the wall, we know that the general solution for the motion is the sum of two functions, $F(x - ct)$ and $G(x + ct)$, the first representing a wave travelling one way in the string, and the second a wave travelling the other way in the string:

$$y = F(x - ct) + G(x + ct) \quad (49.1)$$

is the general solution for any string. But we have next to satisfy the condition that the string does not move at one end. If we put $x = 0$ in Eq. (49.1) and examine y for any value of t , we get $y = F(-ct) + G(+ct)$. Now if this is to be zero for all times, it means that the function $G(ct)$ must be $-F(-ct)$. In other words, G of anything must be $-F$ of minus that same thing. If this result is put back into Eq. (49.1), we find that the solution for the problem is

$$y = F(x - ct) - F(-x - ct). \quad (49.2)$$

It is easy to check that we will get $y = 0$ if we set $x = 0$.

Figure 49-1 shows a wave travelling in the negative x -direction near $x = 0$, and a hypothetical wave travelling in the other direction reversed in sign and on the other side of the origin. We say hypothetical because, of course, there is no string to vibrate on that side of the origin. The total motion of the string is to be regarded as the sum of these two waves in the region of positive x . As they reach the origin, they will always cancel at $x = 0$, and finally the second (reflected) wave will be the only one to exist for positive x and it will, of course, be travelling in the opposite direction. These results are equivalent to the following statement: if a wave reaches the clamped end of a string, it will be reflected with a change in sign. Such a reflection can always be understood by imagining that what is coming to the end of the string comes out upside down from behind the wall. In short, if we assume that the string is infinite and that whenever we have a wave going one way we have another one going the other way with the stated symmetry, the displacement at $x = 0$ will always be zero and it would make no difference if we clamped the string there.

The next point to be discussed is the reflection of a periodic wave. Suppose that the wave represented by $F(x - ct)$ is a sine wave and has been reflected; then the reflected wave $-F(-x - ct)$ is also a sine wave of the same frequency, but travelling in the opposite direction. This situation can be most simply described by using the complex function notation: $F(x - ct) = e^{i\omega(t - x/c)}$ and

49-1 The reflection of waves

49-2 Confined waves, with natural frequencies

49-3 Modes in two dimensions

49-4 Coupled pendulums

49-5 Linear systems

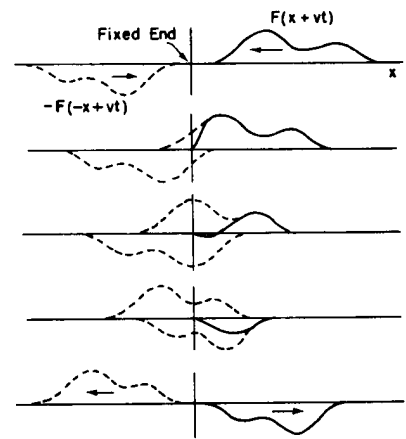


Fig. 49-1. Reflection of a wave at a fixed end as a superposition of two travelling waves.

$F(-x - ct) = e^{i\omega(t+x/c)}$. It can be seen that if these are substituted in (49.2) and if x is set equal to 0, then $y = 0$ for all values of t , so it satisfies the necessary condition. Because of the properties of exponentials, this can be written in a simpler form:

$$y = e^{i\omega t}(e^{-i\omega x/c} - e^{i\omega x/c}) = -2ie^{i\omega t} \sin(\omega x/c). \quad (49.3)$$

There is something interesting and new here, in that this solution tells us that if we look at any fixed x , the string oscillates at frequency ω . No matter where this point is, the frequency is the same! But there are some places, in particular wherever $\sin(\omega x/c) = 0$, where there is no displacement at all. Furthermore, if at any time t we take a snapshot of the vibrating string, the picture will be a sine wave. However, the displacement of this sine wave will depend upon the time t . From inspection of Eq. (49.3) we can see that the length of one cycle of the sine wave is equal to the wavelength of either of the superimposed waves:

$$\lambda = 2\pi c/\omega. \quad (49.4)$$

The points where there is no motion satisfy the condition $\sin(\omega x/c) = 0$, which means that $(\omega x/c) = 0, \pi, 2\pi, \dots, n\pi \dots$. These points are called *nodes*. Between any two successive nodes, every point moves up and down sinusoidally, but the pattern of motion stays fixed in space. This is the fundamental characteristic of what we call a *mode*. If one can find a pattern of motion which has the property that at any point the object moves perfectly sinusoidally, and that all points move at the same frequency (though some will move more than others), then we have what is called a mode.

49-2 Confined waves, with natural frequencies

The next interesting problem is to consider what happens if the string is held at both ends, say at $x = 0$ and $x = L$. We can begin with the idea of the reflection of waves, starting with some kind of a bump moving in one direction. As time goes on, we would expect the bump to get near one end, and as time goes still further it will become a kind of little wobble, because it is combining with the reversed-image bump which is coming from the other side. Finally the original bump will disappear and the image bump will move in the other direction to repeat the process at the other end. This problem has an easy solution, but an interesting question is whether we can have a sinusoidal motion (the solution just described is *periodic*, but of course it is not *sinusoidally* periodic). Let us try to put a sinusoidally periodic wave on a string. If the string is tied at one end, we know it must look like our earlier solution (49.3). If it is tied at the other end, it has to look the same at the other end. So the only possibility for periodic sinusoidal motion is that the sine wave must neatly fit into the string length. If it does not fit into the string length, then it is not a natural frequency at which the string can continue to oscillate. In short, if the string is started with a sine wave shape that just fits in, then it will continue to keep that perfect shape of a sine wave and will oscillate harmonically at some frequency.

Mathematically, we can write $\sin kx$ for the shape, where k is equal to the factor (ω/c) in Eqs. (49.3) and (49.4), and this function will be zero at $x = 0$. However, it must also be zero at the other end. The significance of this is that k is no longer arbitrary, as was the case for the half-open string. With the string closed at both ends, the only possibility is that $\sin(kL) = 0$, because this is the only condition that will keep both ends fixed. Now in order for a sine to be zero, the angle must be either 0, π , 2π , or some other integral multiple of π . The equation

$$kL = n\pi \quad (49.5)$$

will, therefore, give any one of the possible k 's, depending on what integer is put in. For each of the k 's there is a certain frequency ω , which, according to (49.3), is simply

$$\omega = kc = n\pi c/L. \quad (49.6)$$

So we have found the following: that a string has a property that it can have sinusoidal motions, *but only at certain frequencies*. This is the most important characteristic of confined waves. No matter how complicated the system is, it always turns out that there are some patterns of motion which have a perfect sinusoidal time dependence, but with frequencies that are a property of the particular system and the nature of its boundaries. In the case of the string we have many different possible frequencies, each one, by definition, corresponding to a mode, because a mode is a pattern of motion which repeats itself sinusoidally. Figure 49-2 shows the first three modes for a string. For the first mode the wavelength λ is $2L$. This can be seen if one continues the wave out to $x = 2L$ to obtain one complete cycle of the sine wave. The angular frequency ω is $2\pi c$ divided by the wavelength, in general, and in this case, since λ is $2L$, the frequency is $\pi c/L$, which is in agreement with (49.6) with $n = 1$. Let us call the first mode frequency ω_1 . Now the next mode shows two loops with one node in the middle. For this mode the wavelength, then, is simply L . The corresponding value of k is twice as great and the frequency is twice as large; it is $2\omega_1$. For the third mode it is $3\omega_1$, and so on. So all the different frequencies of the string are multiples, 1, 2, 3, 4, and so on, of the lowest frequency ω_1 .

Returning now to the general motion of the string, it turns out that any possible motion can always be analyzed by asserting that more than one mode is operating at the same time. In fact, for general motion an infinite number of modes must be excited at the same time. To get some idea of this, let us illustrate what happens when there are two modes oscillating at the same time: Suppose that we have the first mode oscillating as shown by the sequence of pictures in Fig. 49-3, which illustrates the deflection of the string for equally spaced time intervals extending through half a cycle of the lowest frequency.

Now, at the same time, we suppose that there is an oscillation of the second mode also. Figure 49-3 also shows a sequence of pictures of this mode, which at the start is 90° out of phase with the first mode. This means that at the start it has no displacement, but the two halves of the string have oppositely directed velocities. Now we recall a general principle relating to linear systems: if there are any two solutions, then their sum is also a solution. Therefore a third possible motion of the string would be a displacement obtained by adding the two solutions shown in Fig. 49-3. The result, also shown in the figure, begins to suggest the idea of a bump running back and forth between the ends of the string, although with only two modes we cannot make a very good picture of it; more modes are needed. This result is, in fact, a special case of a great principle for linear systems:

Any motion at all can be analyzed by assuming that it is the sum of the motions of all the different modes, combined with appropriate amplitudes and phases.

The importance of the principle derives from the fact that each mode is very simple—it is nothing but a sinusoidal motion in time. It is true that even the general motion of a string is not really very complicated, but there are other systems, for example the whipping of an airplane wing, in which the motion is much more complicated. Nevertheless, even with an airplane wing, we find there is a certain particular way of twisting which has one frequency and other ways of twisting that have other frequencies. If these modes can be found, then the complete motion can always be analyzed as a superposition of harmonic oscillations (except when the whipping is of such degree that the system can no longer be considered as linear).

49-3 Modes in two dimensions

The next example to be considered is the interesting situation of modes in two dimensions. Up to this point we have talked only about one-dimensional situations—a stretched string or sound waves in a tube. Ultimately we should consider three dimensions, but an easier step will be that to two dimensions. Consider for definiteness a rectangular rubber drumhead which is confined so as to have no displacement anywhere on the rectangular edge, and let the dimensions of the rectangle

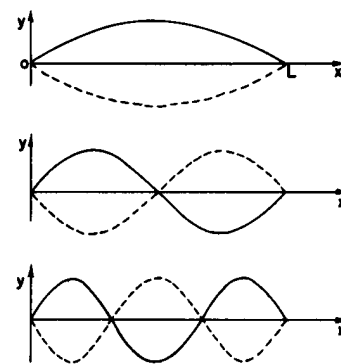


Fig. 49-2. The first three modes of a vibrating string.

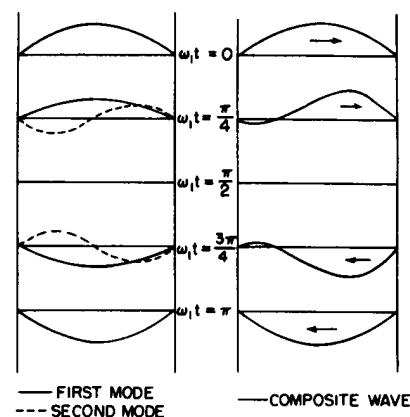


Fig. 49-3. Two modes combine to give a travelling wave.

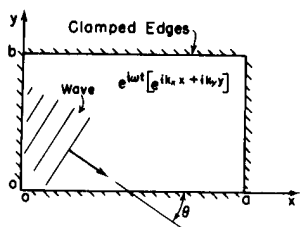


Fig. 49-4. Vibrating rectangular plate.

be a and b , as shown in Fig. 49-4. Now the question is, what are the characteristics of the possible motion? We can start with the same procedure used for the string. If we had no confinement at all, we would expect waves travelling along with some kind of wave motion. For example, $(e^{i\omega t})(e^{-ik_x x + ik_y y})$ would represent a sine wave travelling in some direction which depends on the relative values of k_x and k_y . Now how can we make the x -axis, that is, the line $y = 0$, a node? Using the ideas developed for the one-dimensional string, we can imagine another wave represented by the complex function $(-e^{i\omega t})(e^{-ik_x x - ik_y y})$. The superposition of these waves will give zero displacement at $y = 0$ regardless of the values of x and t . (Although these functions are defined for negative y where there is no drumhead to vibrate, this can be ignored, since the displacement is truly zero at $y = 0$.) In this case we can look upon the second function as the reflected wave.

However, we want a nodal line at $y = b$ as well as at $y = 0$. How do we do that? The solution is related to something we did when studying reflection from crystals. These waves which cancel each other at $y = 0$ will do the same at $y = b$ only if $2b \sin \theta$ is an integral multiple of λ , where θ is the angle shown in Fig. 49-4:

$$m\lambda = 2b \sin \theta, \quad m = 0, 1, 2, \dots \quad (49.7)$$

Now in the same way we can make the y -axis a nodal line by adding two more functions $-(e^{i\omega t})(e^{+ik_x x + ik_y y})$ and $+(e^{i\omega t})(e^{+ik_x x - ik_y y})$, each representing a reflection of one of the other two waves from the $x = 0$ line. The condition for a nodal line at $x = a$ is similar to the one for $y = b$. It is that $2a \cos \theta$ must also be an integral multiple of λ :

$$n\lambda = 2a \cos \theta. \quad (49.8)$$

Then the final result is that the waves bouncing about in the box produce a standing-wave pattern, that is, a definite mode.

So we must satisfy the above two conditions if we are to have a mode. Let us first find the wavelength. This can be obtained by eliminating the angle θ from (49.7) and (49.8) to obtain the wavelength in terms of a , b , n and m . The easiest way to do that is to divide both sides of the respective equations by $2b$ and $2a$, square them, and add the two equations together. The result is $\sin^2 \theta + \cos^2 \theta = 1 = (n\lambda/2a)^2 + (m\lambda/2b)^2$, which can be solved for λ :

$$\frac{1}{\lambda^2} = \frac{n^2}{4a^2} + \frac{m^2}{4b^2}. \quad (49.9)$$

In this way we have determined the wavelength in terms of two integers, and from the wavelength we immediately get the frequency ω , because, as we know, the frequency is equal to $2\pi c$ divided by the wavelength.

This result is interesting and important enough that we should deduce it by a purely mathematical analysis instead of by an argument about the reflections. Let us represent the vibration by a superposition of four waves chosen so that the four lines $x = 0$, $x = a$, $y = 0$, and $y = b$ are all nodes. In addition we shall require that all waves have the same frequency, so that the resulting motion will represent a mode. From our earlier treatment of light reflection we know that $(e^{i\omega t})(e^{-ik_x x + ik_y y})$ represents a wave travelling in the direction indicated in Fig. 49-4. Equation (49.6), that is, $k = \omega/c$, still holds, provided

$$k^2 = k_x^2 + k_y^2. \quad (49.10)$$

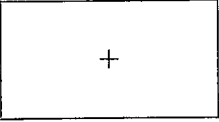
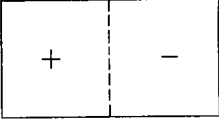
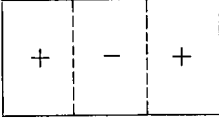
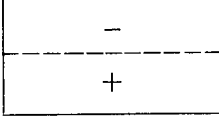
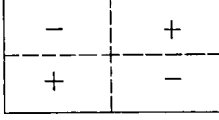
It is clear from the figure that $k_x = k \cos \theta$ and $k_y = k \sin \theta$.

Now our equation for the displacement, say ϕ , of the rectangular drumhead takes on the grand form

$$\phi = [e^{i\omega t}][e^{(-ik_x x + ik_y y)} - e^{(+ik_x x + ik_y y)} - e^{(-ik_x x - ik_y y)} + e^{(+ik_x x - ik_y y)}]. \quad (49.11a)$$

Although this looks rather a mess, the sum of these things now is not very hard.

Table 49-1

Mode shape	m	n	$(\omega/\omega_0)^2$	ω/ω_0
	1	1	1.25	1.12
	1	2	2.00	1.41
	1	3	3.25	1.80
	2	1	4.25	2.06
	2	2	5.00	2.24

The exponentials can be combined to give sine functions, so that the displacement turns out to be

$$\phi = [-4 \sin k_x x \sin k_y y][e^{i\omega t}]. \quad (49.11b)$$

In other words, it is a sinusoidal oscillation, all right, with a pattern that is also sinusoidal in both the x - and the y -direction. Our boundary conditions are of course satisfied at $x = 0$ and $y = 0$. We also want ϕ to be zero when $x = a$ and when $y = b$. Therefore we have to put in two other conditions: $k_x a$ must be an integral multiple of π , and $k_y b$ must be another integral multiple of π . Since we have seen that $k_x = k \cos \theta$ and $k_y = k \sin \theta$, we immediately get equations (49.7) and (49.8) and from these the final result (49.9).

Now let us take as an example a rectangle whose width is twice the height. If we take $a = 2b$ and use Eqs. (49.4) and (49.9), we can calculate the frequencies of all of the modes:

$$\omega^2 = \left(\frac{\pi c}{b}\right)^2 \frac{4m^2 + n^2}{4}. \quad (49.12)$$

Table 49-1 lists a few of the simple modes and also shows their shape in a qualitative way.

The most important point to be emphasized about this particular case is that the frequencies are not multiples of each other, nor are they multiples of any number. The idea that the natural frequencies are harmonically related is not generally true. It is not true for a system of more than one dimension, nor is it true for one-dimensional systems which are more complicated than a string with uniform density and tension. A simple example of the latter is a hanging chain in which the tension is higher at the top than at the bottom. If such a chain is set in harmonic oscillation, there are various modes and frequencies, but the frequencies are not simple multiples of any number, nor are the mode shapes sinusoidal.

The modes of more complicated systems are still more elaborate. For example, inside the mouth we have a cavity above the vocal cords, and by moving the tongue and the lips, and so forth, we make an open-ended pipe or a closed-ended pipe of different diameters and shapes; it is a terribly complicated resonator, but

it is a resonator nevertheless. Now when one talks with the vocal cords, they are made to produce some kind of tone. The tone is rather complicated and there are many sounds coming out, but the cavity of the mouth further modifies that tone because of the various resonant frequencies of the cavity. For instance, a singer can sing various vowels, a, or o, or oo, and so forth, at the same pitch, but they sound different because the various harmonics are in resonance in this cavity to different degrees. The very great importance of the resonant frequencies of a cavity in modifying the voice sounds can be demonstrated by a simple experiment. Since the speed of sound goes as the reciprocal of the square root of the density, the speed of sound may be varied by using different gases. If one uses helium instead of air, so that the density is lower, the speed of sound is much higher, and all the frequencies of a cavity will be raised. Consequently if one fills one's lungs with helium before speaking, the character of his voice will be drastically altered even though the vocal cords may still be vibrating at the same frequency.

49-4 Coupled pendulums

Finally we should emphasize that not only do modes exist for complicated continuous systems, but also for very simple mechanical systems. A good example is the system of two coupled pendulums discussed in the preceding chapter. In that chapter it was shown that the motion could be analyzed as a superposition of two harmonic motions with different frequencies. So even this system can be analyzed in terms of harmonic motions or modes. The string has an infinite number of modes and the two-dimensional surface also has an infinite number of modes. In a sense it is a double infinity, if we know how to count infinities. But a simple mechanical thing which has only two degrees of freedom, and requires only two variables to describe it, has only two modes.

Let us make a mathematical analysis of these two modes for the case where the pendulums are of equal length. Let the displacement of one be x , and the displacement of the other be y , as shown in Fig. 49-5. Without a spring, the force on the first mass is proportional to the displacement of that mass, because of gravity. There would be, if there were no spring, a certain natural frequency ω_0 for this one alone. The equation of motion without a spring would be

$$m \frac{d^2 x}{dt^2} = -m\omega_0^2 x. \quad (49.13)$$

The other pendulum would swing in the same way if there were no spring. In addition to the force of restoration due to gravitation, there is an additional force pulling the first mass. That force depends upon the excess distance of x over y and is proportional to that difference, so it is some constant which depends on the geometry, times $(x - y)$. The same force in reverse sense acts on the second mass. The equations of motion that have to be solved are therefore

$$m \frac{d^2 x}{dt^2} = -m\omega_0^2 x - k(x - y), \quad m \frac{d^2 y}{dt^2} = -m\omega_0^2 y - k(y - x). \quad (49.14)$$

In order to find a motion in which both of the masses move at the same frequency, we must determine how much each mass moves. In other words, pendulum x and pendulum y will oscillate at the same frequency, but their amplitudes must have certain values, A and B , whose relation is fixed. Let us try this solution:

$$x = Ae^{i\omega t}, \quad y = Be^{i\omega t}. \quad (49.15)$$

If these are substituted in Eqs. (49.14) and similar terms are collected, the results are

$$\begin{aligned} \left(\omega^2 - \omega_0^2 - \frac{k}{m} \right) A &= -\frac{k}{m} B, \\ \left(\omega^2 - \omega_0^2 - \frac{k}{m} \right) B &= -\frac{k}{m} A. \end{aligned} \quad (49.16)$$

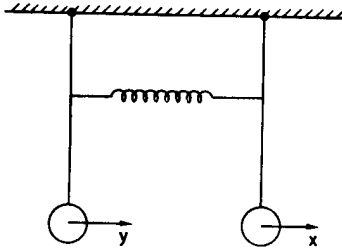


Fig. 49-5. Two coupled pendulums.

The equations as written have had the common factor $e^{i\omega t}$ removed and have been divided by m .

Now we see that we have two equations for what looks like two unknowns. But there really are not *two* unknowns, because the whole size of the motion is something that we cannot determine from these equations. The above equations can determine only the *ratio* of A to B , *but they must both give the same ratio*. The necessity for both of these equations to be consistent is a requirement that the frequency be something very special.

In this particular case this can be worked out rather easily. If the two equations are multiplied together, the result is

$$\left(\omega^2 - \omega_0^2 - \frac{k}{m}\right)^2 AB = \left(\frac{k}{m}\right)^2 AB. \quad (49.17)$$

The term AB can be removed from both sides unless A and B are zero, which means there is no motion at all. If there is motion, then the other terms must be equal, giving a quadratic equation to solve. The result is that there are two possible frequencies:

$$\omega_1^2 = \omega_0^2, \quad \omega_2^2 = \omega_0^2 + \frac{2k}{m}. \quad (49.18)$$

Furthermore, if these values of frequency are substituted back into Eq. (49.16), we find that for the first frequency $A = B$, and for the second frequency $A = -B$. These are the “mode shapes,” as can be readily verified by experiment.

It is clear that in the first mode, where $A = B$, the spring is never stretched, and both masses oscillate at the frequency ω_0 , as though the spring were absent. In the other solution, where $A = -B$, the spring contributes a restoring force and raises the frequency. A more interesting case results if the pendulums have different lengths. The analysis is very similar to that given above, and is left as an exercise for the reader.

49-5 Linear systems

Now let us summarize the ideas discussed above, which are all aspects of what is probably the most general and wonderful principle of mathematical physics. If we have a linear system whose character is independent of the time, then the motion does not have to have any particular simplicity, and in fact may be exceedingly complex, but there are very special motions, usually a series of special motions, in which the whole pattern of motion varies exponentially with the time. For the vibrating systems that we are talking about now, the exponential is imaginary, and instead of saying “exponentially” we might prefer to say “sinusoidally” with time. However, one can be more general and say that the motions will vary exponentially with the time in very special modes, with very special shapes. The most general motion of the system can always be represented as a superposition of motions involving each of the different exponentials.

This is worth stating again for the case of sinusoidal motion: a linear system need not be moving in a purely sinusoidal motion, i.e., at a definite single frequency, but no matter how it does move, this motion can be represented as a superposition of pure sinusoidal motions. The frequency of each of these motions is a characteristic of the system, and the pattern or waveform of each motion is also a characteristic of the system. The general motion in any such system can be characterized by giving the strength and the phase of each of these modes, and adding them all together. Another way of saying this is that any linear vibrating system is equivalent to a set of independent harmonic oscillators, with the natural frequencies corresponding to the modes.

We conclude this chapter by remarking on the connection of modes with quantum mechanics. In quantum mechanics the vibrating object, or the thing that varies in space, is the amplitude of a probability function that gives the probability of finding an electron, or system of electrons, in a given configuration. This amplitude function can vary in space and time, and satisfies, in fact, a linear equation.

But in quantum mechanics there is a transformation, in that what we call frequency of the probability amplitude is equal, in the classical idea, to energy. Therefore we can translate the principle stated above to this case by taking the word *frequency* and replacing it with *energy*. It becomes something like this: a quantum-mechanical system, for example an atom, need not have a definite energy, just as a simple mechanical system does not have to have a definite frequency; but no matter how the system behaves, its behavior can always be represented as a superposition of states of definite energy. The energy of each state is a characteristic of the atom, and so is the pattern of amplitude which determines the probability of finding particles in different places. The general motion can be described by giving the amplitude of each of these different energy states. This is the origin of energy levels in quantum mechanics. Since quantum mechanics is represented by waves, in the circumstance in which the electron does not have enough energy to ultimately escape from the proton, they are *confined waves*. Like the confined waves of a string, there are definite frequencies for the solution of the wave equation for quantum mechanics. The quantum-mechanical interpretation is that these are definite *energies*. Therefore a quantum-mechanical system, because it is represented by waves, can have definite states of fixed energy; examples are the energy levels of various atoms.

Harmonics

50-1 Musical tones

Pythagoras is said to have discovered the fact that two similar strings under the same tension and differing only in length, when sounded together give an effect that is pleasant to the ear *if* the lengths of the strings are in the ratio of two small integers. If the lengths are as one is to two, they then correspond to the octave in music. If the lengths are as two is to three, they correspond to the interval between *C* and *G*, which is called a fifth. These intervals are generally accepted as “pleasant” sounding chords.

Pythagoras was so impressed by this discovery that he made it the basis of a school—Pythagoreans they were called—which held mystic beliefs in the great powers of numbers. It was believed that something similar would be found out about the planets—or “spheres.” We sometimes hear the expression: “the music of the spheres.” The idea was that there would be some numerical relationships between the orbits of the planets or between other things in nature. People usually think that this is just a kind of superstition held by the Greeks. But is it so different from our own scientific interest in quantitative relationships? Pythagoras’ discovery was the first example, outside geometry, of any numerical relationship in nature. It must have been very surprising to suddenly discover that there was a *fact* of nature that involved a simple numerical relationship. Simple measurements of lengths gave a prediction about something which had no apparent connection to geometry—the production of pleasant sounds. This discovery led to the extension that perhaps a good tool for understanding nature would be arithmetic and mathematical analysis. The results of modern science justify that point of view.

Pythagoras could only have made his discovery by making an experimental observation. Yet this important aspect does not seem to have impressed him. If it had, physics might have had a much earlier start. (It is always easy to look back at what someone else has done and to decide what he *should* have done!)

We might remark on a third aspect of this very interesting discovery: that the discovery had to do with two notes that *sound pleasant* to the ear. We may question whether we are any better off than Pythagoras in understanding *why* only certain sounds are pleasant to our ear. The general theory of aesthetics is probably no further advanced now than in the time of Pythagoras. In this one discovery of the Greeks, there are the three aspects: experiment, mathematical relationships, and aesthetics. Physics has made great progress on only the first two parts. This chapter will deal with our present-day understanding of the discovery of Pythagoras.

Among the sounds that we hear, there is one kind that we call *noise*. Noise corresponds to a sort of irregular vibration of the eardrum that is produced by the irregular vibration of some object in the neighborhood. If we make a diagram to indicate the pressure of the air on the eardrum (and, therefore, the displacement of the drum) as a function of time, the graph which corresponds to a noise might look like that shown in Fig. 50-1(a). (Such a noise might correspond roughly to the sound of a stamped foot.) The sound of *music* has a different character. Music is characterized by the presence of more-or-less *sustained tones*—or musical “notes.” (Musical instruments may make noises as well!) The tone may last for a relatively short time, as when a key is pressed on a piano, or it may be sustained almost indefinitely, as when a flute player holds a long note.

What is the special character of a musical note from the point of view of the pressure in the air? A musical note differs from a noise in that there is a periodicity in its graph. There is some uneven shape to the variation of the air pressure with

50-1 Musical tones

50-2 The Fourier series

50-3 Quality and consonance

50-4 The Fourier coefficients

50-5 The energy theorem

50-6 Nonlinear responses

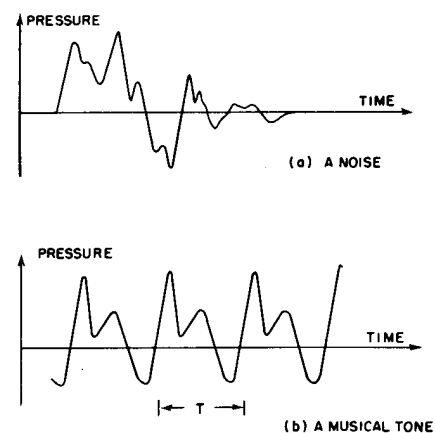


Fig. 50-1. Pressure as a function of time for (a) a noise, and (b) a musical tone.

time, and the shape repeats itself over and over again. An example of a pressure-time function that would correspond to a musical note is shown in Fig. 50-1(b).

Musicians will usually speak of a musical tone in terms of three characteristics: the loudness, the pitch, and the "quality." The "loudness" is found to correspond to the magnitude of the pressure changes. The "pitch" corresponds to the period of time for one repetition of the basic pressure function. ("Low" notes have longer periods than "high" notes.) The "quality" of a tone has to do with the differences we may still be able to hear between two notes of the same loudness and pitch. An oboe, a violin, or a soprano are still distinguishable even when they sound notes of the same pitch. The quality has to do with the structure of the repeating pattern.

Let us consider, for a moment, the sound produced by a vibrating string. If we pluck the string, by pulling it to one side and releasing it, the subsequent motion will be determined by the motions of the waves we have produced. We know that these waves will travel in both directions, and will be reflected at the ends. They will slosh back and forth for a long time. No matter how complicated the wave is, however, it will repeat itself. The period of repetition is just the time T required for the wave to travel two full lengths of the string. For that is just the time required for any wave, once started, to reflect off each end and return to its starting position, and be proceeding in the original direction. The time is the same for waves which start out in either direction. Each point on the string will, then, return to its starting position after one period, and again one period later, etc. The sound wave produced must also have the same repetition. We see why a plucked string produces a musical tone.

50-2 The Fourier series

We have discussed in the preceding chapter another way of looking at the motion of a vibrating system. We have seen that a string has various natural modes of oscillation, and that any particular kind of vibration that may be set up by the starting conditions can be thought of as a combination—in suitable proportions—of several of the natural modes, oscillating together. For a string we found that the normal modes of oscillation had the frequencies ω_0 , $2\omega_0$, $3\omega_0$, . . . The most general motion of a plucked string, therefore, is composed of the sum of a sinusoidal oscillation at the fundamental frequency ω_0 , another at the second harmonic frequency $2\omega_0$, another at the third harmonic $3\omega_0$, etc. Now the fundamental mode repeats itself every period $T_1 = 2\pi/\omega_0$. The second harmonic mode repeats itself every $T_2 = 2\pi/2\omega_0$. It *also* repeats itself every $T_1 = 2T_2$, after *two* of its periods. Similarly, the third harmonic mode repeats itself after a time T_1 which is 3 of its periods. We see again why a plucked string repeats its whole pattern with a periodicity of T_1 . It produces a musical tone.

We have been talking about the motion of the string. But the *sound*, which is the motion of the air, is produced by the motion of the string, so its vibrations too must be composed of the same harmonics—though we are no longer thinking about the normal modes of the air. Also, the relative strength of the harmonics may be different in the air than in the string, particularly if the string is "coupled" to the air via a sounding board. The efficiency of the coupling to the air is different for different harmonics.

If we let $f(t)$ represent the air pressure as a function of time for a musical tone [such as that in Fig. 50-1(b)], then we expect that $f(t)$ can be written as the sum of a number of simple harmonic functions of time—like $\cos \omega t$ —for each of the various harmonic frequencies. If the period of the vibration is T , the fundamental angular frequency will be $\omega = 2\pi/T$, and the harmonics will be 2ω , 3ω , etc.

There is one slight complication. For each frequency we may expect that the starting phases will not necessarily be the same for all frequencies. We should, therefore, use functions like $\cos(\omega t + \phi)$. It is, however, simpler to use instead both the sine and cosine functions for *each* frequency. We recall that

$$\cos(\omega t + \phi) = (\cos \phi \cos \omega t - \sin \phi \sin \omega t) \quad (50.1)$$

and since ϕ is a constant, *any* sinusoidal oscillation at the frequency ω can be written as the sum of a term with $\cos \omega t$ and another term with $\sin \omega t$.

We conclude, then, that *any* function $f(t)$ that is periodic with the period T can be written mathematically as

$$\begin{aligned} f(t) = & a_0 \\ & + a_1 \cos \omega t + b_1 \sin \omega t \\ & + a_2 \cos 2\omega t + b_2 \sin 2\omega t \\ & + a_3 \cos 3\omega t + b_3 \sin 3\omega t \\ & + \dots + \dots \end{aligned} \quad (50.2)$$

where $\omega = 2\pi/T$ and the a 's and b 's are numerical constants which tell us how much of each component oscillation is present in the oscillation $f(t)$. We have added the "zero-frequency" term a_0 so that our formula will be completely general, although it is usually zero for a musical tone. It represents a shift of the average value (that is, the "zero" level) of the sound pressure. With it our formula can take care of any case. The equality of Eq. (50.2) is represented schematically in Fig. 50-2. (The amplitudes, a_n and b_n , of the harmonic functions must be suitably chosen. They are shown schematically and without any particular scale in the figure.) The series (50.2) is called the *Fourier series* for $f(t)$.

We have said that *any* periodic function can be made up in this way. We should correct that and say that any sound wave, or any function we ordinarily encounter in physics, can be made up of such a sum. The mathematicians can invent functions which cannot be made up of simple harmonic functions—for instance, a function that has a "reverse twist" so that it has two values for some values of t ! We need not worry about such functions here.

50-3 Quality and consonance

Now we are able to describe what it is that determines the "quality" of a musical tone. It is the relative amounts of the various harmonics—the values of the a 's and b 's. A tone with only the first harmonic is a "pure" tone. A tone with many strong harmonics is a "rich" tone. A violin produces a different proportion of harmonics than does an oboe.

We can "manufacture" various musical tones if we connect several "oscillators" to a loudspeaker. (An oscillator usually produces a nearly pure simple harmonic function.) We should choose the frequencies of the oscillators to be ω , 2ω , 3ω , etc. Then by adjusting the volume control on each oscillator, we can add in any amount we wish of each harmonic—thereby producing tones of different quality. An electric organ works in much this way. The "keys" select the frequency of the fundamental oscillator and the "stops" are switches that control the relative proportions of the harmonics. By throwing these switches, the organ can be made to sound like a flute, or an oboe, or a violin.

It is interesting that to produce such "artificial" tones we need only one oscillator for each frequency—we do not need separate oscillators for the sine and cosine components. The ear is not very sensitive to the relative phases of the harmonics. It pays attention mainly to the *total* of the sine and cosine parts of each frequency. Our analysis is more accurate than is necessary to explain the *subjective* aspect of music. The response of a microphone or other physical instrument does depend on the phases, however, and our complete analysis may be needed to treat such cases.

The "quality" of a spoken sound also determines the vowel sounds that we recognize in speech. The shape of the mouth determines the frequencies of the natural modes of vibration of the air in the mouth. Some of these modes are set into vibration by the sound waves from the vocal chords. In this way, the amplitudes of some of the harmonics of the sound are increased with respect to others. When we change the shape of our mouth, harmonics of different frequencies are given preference. These effects account for the difference between an "e-e-e" sound and an "a-a-a" sound.

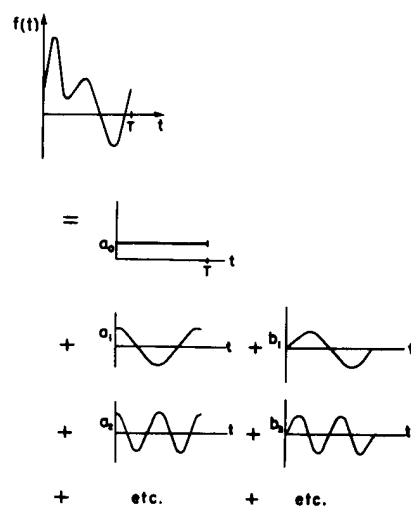


Fig. 50-2. Any periodic function $f(t)$ is equal to a sum of simple harmonic functions.

We all know that a particular vowel sound—say “e-e-e”—still “sounds like” the same vowel whether we say (or sing) it at a high or a low pitch. From the mechanism we describe, we would expect that *particular* frequencies are emphasized when we shape our mouth for an “e-e-e,” and that they *do not* change as we change the pitch of our voice. So the relation of the important harmonics to the fundamental—that is, the “quality”—changes as we change pitch. Apparently the mechanism by which we recognize speech is not based on specific harmonic relationships.

What should we say now about Pythagoras’ discovery? We understand that two similar strings with lengths in the ratio of 2 to 3 will have fundamental frequencies in the ratio 3 to 2. But why should they “sound pleasant” together? Perhaps we should take our clue from the frequencies of the harmonics. The second harmonic of the lower shorter string will have the *same* frequency as the third harmonic of the longer string. (It is easy to show—or to believe—that a plucked string produces strongly the several lowest harmonics.)

Perhaps we should make the following rules. Notes sound consonant when they have harmonics with the same frequency. Notes sound dissonant if their upper harmonics have frequencies near to each other but far enough apart that there are rapid beats between the two. Why beats do not sound pleasant, and why unison of the upper harmonics does sound pleasant, is something that we do not know how to define or describe. We cannot say from this knowledge of what *sounds* good, what ought, for example, to *smell* good. In other words, our understanding of it is not anything more general than the statement that when they are in unison they sound good. It does not permit us to deduce anything more than the properties of concordance in music.

It is easy to check on the harmonic relationships we have described by some simple experiments with a piano. Let us label the 3 successive C’s near the middle of the keyboard by C, C’, and C’’, and the G’s just above by G, G’, and G’’. Then the fundamentals will have relative frequencies as follows:

$$\begin{array}{ll} \text{C} - 2 & \text{G} - 3 \\ \text{C}' - 4 & \text{G}' - 6 \\ \text{C}'' - 8 & \text{G}'' - 12 \end{array}$$

These harmonic relationships can be demonstrated in the following way: Suppose we press C’ *slowly*—so that it does not sound but we cause the damper to be lifted. If we then sound C, it will produce its own fundamental *and* some second harmonic. The second harmonic will set the strings of C’ into vibration. If we now release C (keeping C’ pressed) the damper will stop the vibration of the C strings, and we can hear (softly) the note C’ as it dies away. In a similar way, the third harmonic of C can cause a vibration of G’. Or the sixth of C (now getting much weaker) can set up a vibration in the fundamental of G’.

A somewhat different result is obtained if we press G quietly and then sound C’. The third harmonic of C’ will correspond to the fourth harmonic of G, so *only* the fourth harmonic of G will be excited. We can hear (if we listen closely) the sound of G’’, which is two octaves above the G we have pressed! It is easy to think up many more combinations for this game.

We may remark in passing that the major scale can be defined just by the condition that the three major chords (F-A-C); (C-E-G); and (G-B-D) *each* represent tone sequences with the frequency ratio (4: 5: 6). These ratios—plus the fact that an octave (C-C’, B-B’, etc.) has the ratio 1: 2—determine the whole scale for the “ideal” case, or for what is called “just intonation.” Keyboard instruments like the piano are *not* usually tuned in this manner, but a little “fudging” is done so that the frequencies are *approximately* correct for all possible starting tones. For this tuning, which is called “tempered,” the octave (still 1: 2) is divided into 12 equal intervals for which the frequency ratio is $(2)^{1/12}$. A fifth no longer has the frequency ratio 3/2, but $2^{7/12} = 1.499$, which is apparently close enough for most ears.

We have stated a rule for consonance in terms of the coincidence of harmonics. Is this coincidence perhaps the *reason* that two notes are consonant? One worker has claimed that two *pure* tones—tones carefully manufactured to be free of harmonics—do not give the *sensations* of consonance or dissonance as the relative frequencies are placed at or near the expected ratios. (Such experiments are difficult because it is difficult to manufacture pure tones, for reasons that we shall see later.) We cannot still be certain whether the ear is matching harmonics or doing arithmetic when we decide that we like a sound.

50-4 The Fourier coefficients

Let us return now to the idea that any note—that is, a *periodic* sound—can be represented by a suitable combination of harmonics. We would like to show how we can find out what amount of each harmonic is required. It is, of course, easy to compute $f(t)$, using Eq. (50.2), if we are *given* all the coefficients a and b . The question now is, if we are given $f(t)$ how can we know what the coefficients of the various harmonic terms should be? (It is easy to make a cake from a recipe; but can we write down the recipe if we are given a cake?)

Fourier discovered that it was not really very difficult. The term a_0 is certainly easy. We have already said that it is just the average value of $f(t)$ over one period (from $t = 0$ to $t = T$). We can easily see that this is indeed so. The average value of a sine or cosine function over one period is zero. Over two, or three, or any whole number of periods, it is also zero. So the average value of all of the terms on the right-hand side of Eq. (50.2) is zero, except for a_0 . (Recall that we must choose $\omega = 2\pi/T$.)

Now the average of a sum is the sum of the averages. So the average of $f(t)$ is just the average of a_0 . But a_0 is a *constant*, so its average is just the same as its value. Recalling the definition of an average, we have

$$a_0 = \frac{1}{T} \int_0^T f(t) dt. \quad (50.3)$$

The other coefficients are only a little more difficult. To find them we can use a trick discovered by Fourier. Suppose we multiply both sides of Eq. (50.2) by some harmonic function—say by $\cos 7\omega t$. We have then

$$\begin{aligned} f(t) \cdot \cos 7\omega t &= a_0 \cdot \cos 7\omega t \\ &+ a_1 \cos \omega t \cdot \cos 7\omega t + b_1 \sin \omega t \cdot \cos 7\omega t \\ &+ a_2 \cos 2\omega t \cdot \cos 7\omega t + b_2 \sin 2\omega t \cdot \cos 7\omega t \\ &+ \cdots \qquad \qquad \qquad + \cdots \\ &+ a_7 \cos 7\omega t \cdot \cos 7\omega t + b_7 \sin 7\omega t \cdot \cos 7\omega t \\ &+ \cdots \qquad \qquad \qquad + \cdots \end{aligned} \quad (50.4)$$

Now let us average both sides. The average of $a_0 \cos 7\omega t$ over the time T is proportional to the average of a cosine over 7 whole periods. But that is just zero. The average of *almost all* of the rest of the terms is *also* zero. Let us look at the a_1 term. We know, in general, that

$$\cos A \cos B = \frac{1}{2} \cos (A + B) + \frac{1}{2} \cos (A - B). \quad (50.5)$$

The a_1 term becomes

$$\frac{1}{2} a_1 (\cos 8\omega t + \cos 6\omega t). \quad (50.6)$$

We thus have two cosine terms, one with 8 full periods in T and the other with 6. *They both average to zero.* The average of the a_1 term is therefore zero.

For the a_2 term, we would find $a_2 \cos 9\omega t$ and $a_2 \cos 5\omega t$, each of which also averages to zero. For the a_9 term, we would find $\cos 16\omega t$ and $\cos (-2\omega t)$. But $\cos (-2\omega t)$ is the same as $\cos 2\omega t$, so both of these have zero averages. It is clear

that *all* of the a terms will have a zero average *except* one. And that one is the a_7 term. For this one we have

$$\frac{1}{2}a_7(\cos 14\omega t + \cos 0). \quad (50.7)$$

The cosine of zero is one, and its average, of course, is one. So we have the result that the average of all of the a terms of Eq. (50.4) equals $\frac{1}{2}a_7$.

The b terms are even easier. When we multiply by any cosine term like $\cos n\omega t$, we can show by the same method that *all* of the b terms have the average value zero.

We see that Fourier's "trick" has acted like a sieve. When we multiply by $\cos 7\omega t$ and average, all terms drop out except a_7 , and we find that

$$\text{Average } [f(t) \cdot \cos 7\omega t] = a_7/2, \quad (50.8)$$

or

$$a_7 = \frac{2}{T} \int_0^T f(t) \cdot \cos 7\omega t \, dt. \quad (50.9)$$

We shall leave it for the reader to show that the coefficient b_7 can be obtained by multiplying Eq. (50.2) by $\sin 7\omega t$ and averaging both sides. The result is

$$b_7 = \frac{2}{T} \int_0^T f(t) \cdot \sin 7\omega t \, dt. \quad (50.10)$$

Now what is true for 7 we expect is true for any integer. So we can summarize our proof and result in the following more elegant mathematical form. If m and n are integers other than zero, and if $\omega = 2\pi/T$, then

$$\text{I. } \int_0^T \sin n\omega t \cos m\omega t \, dt = 0. \quad (50.11)$$

$$\text{II. } \int_0^T \cos n\omega t \cos m\omega t \, dt = \begin{cases} 0 & \text{if } n \neq m. \\ T/2 & \text{if } n = m. \end{cases} \quad (50.12)$$

$$\text{III. } \int_0^T \sin n\omega t \sin m\omega t \, dt = \begin{cases} 0 & \text{if } n \neq m. \\ T/2 & \text{if } n = m. \end{cases}$$

$$\text{IV. } f(t) = a_0 + \sum_{n=1}^{\infty} a_n \cos n\omega t + \sum_{n=1}^{\infty} b_n \sin n\omega t. \quad (50.13)$$

$$\text{V. } a_0 = \frac{1}{T} \int_0^T f(t) \cdot dt. \quad (50.14)$$

$$a_n = \frac{2}{T} \int_0^T f(t) \cdot \cos n\omega t \, dt. \quad (50.15)$$

$$b_n = \frac{2}{T} \int_0^T f(t) \cdot \sin n\omega t \, dt. \quad (50.16)$$

In earlier chapters it was convenient to use the exponential notation for representing simple harmonic motion. Instead of $\cos \omega t$ we used $\text{Re } e^{i\omega t}$, the real part of the exponential function. We have used cosine and sine functions in this chapter because it made the derivations perhaps a little clearer. Our final result of Eq. (50.13) can, however, be written in the compact form

$$f(t) = \text{Re } \sum_{n=0}^{\infty} \hat{a}_n e^{in\omega t}, \quad (50.17)$$

where \hat{a}_n is the complex number $a_n - ib_n$ (with $b_0 = 0$). If we wish to use the same notation throughout, we can write also

$$a_n = \frac{2}{T} \int_0^T f(t) e^{-in\omega t} \, dt \quad (n \geq 1). \quad (50.18)$$

We now know how to “analyze” a periodic wave into its harmonic components. The procedure is called *Fourier analysis*, and the separate terms are called Fourier components. We have *not* shown, however, that once we find all of the Fourier components and add them together, we do indeed get back our $f(t)$. The mathematicians have shown, for a wide class of functions, in fact for all that are of interest to physicists, that if we can do the integrals we will get back $f(t)$. There is one minor exception. If the function $f(t)$ is discontinuous, i.e., if it jumps suddenly from one value to another, the Fourier sum will give a value at the breakpoint halfway between the upper and lower values at the discontinuity. So if we have the strange function $f(t) = 0, 0 \leq t < t_0$, and $f(t) = 1$ for $t_0 \leq t \leq T$, the Fourier sum will give the right value everywhere *except* at t_0 , where it will have the value $\frac{1}{2}$ instead of 1. It is rather unphysical anyway to insist that a function should be zero *up to* t_0 , but 1 *right at* t_0 . So perhaps we should make the “rule” for physicists that any discontinuous function (which can only be a simplification of a *real* physical function) should be defined with halfway values at the discontinuities. Then any such function—with any finite number of such jumps—as well as all other physically interesting functions, are given correctly by the Fourier sum.

As an exercise, we suggest that the reader determine the Fourier series for the function shown in Fig. 50-3. Since the function cannot be written in an explicit algebraic form, you will not be able to do the integrals from zero to T in the usual way. The integrals are easy, however, if we separate them into two parts: the integral from zero to $T/2$ (over which $f(t) = 1$) and the integral from $T/2$ to T (over which $f(t) = -1$). The result should be

$$f(t) = \frac{4}{\pi} \left(\sin \omega t + \frac{1}{3} \sin 3\omega t + \frac{1}{5} \sin 5\omega t + \dots \right), \quad (50.19)$$

where $\omega = 2\pi/T$. We thus find that our square wave (with the particular phase chosen) has only odd harmonics, and their amplitudes are in inverse proportion to their frequencies.

Let us check that Eq. (50.19) does indeed give us back $f(t)$ for some value of t . Let us choose $t = T/4$, or $\omega t = \pi/2$. We have

$$f(t) = \frac{4}{\pi} \left(\sin \frac{\pi}{2} + \frac{1}{3} \sin \frac{3\pi}{2} + \frac{1}{5} \sin \frac{5\pi}{2} + \dots \right) \quad (50.20)$$

$$= \frac{4}{\pi} \left(1 - \frac{1}{3} + \frac{1}{5} - \frac{1}{7} + \dots \right) \quad (50.21)$$

The series* has the value $\pi/4$, and we find that $f(t) = 1$.

50-5 The energy theorem

The energy in a wave is proportional to the square of its amplitude. For a wave of complex shape, the energy in one period will be proportional to $\int_0^T f^2(t) dt$. We can also relate this energy to the Fourier coefficients. We write

$$\int_0^T f^2(t) dt = \int_0^T \left[a_0 + \sum_{n=1}^{\infty} a_n \cos n\omega t + \sum_{n=1}^{\infty} b_n \sin n\omega t \right]^2 dt. \quad (50.22)$$

When we expand the square of the bracketed term we will get all possible cross terms, such as $a_5 \cos 5\omega t \cdot b_7 \cos 7\omega t$. We have shown above, however, [Eqs. (50.11) and (50.12)] that the integrals of all such terms over one period is zero.

* The series can be evaluated in the following way. First we remark that $\int_0^x [dx/(1+x^2)] = \tan^{-1} x$. Second, we expand the integrand in a series $1/(1+x^2) = 1 - x^2 + x^4 - x^6 + \dots$. We integrate the series term by term (from zero to x) to obtain $\tan^{-1} x = x - x^3/3 + x^5/5 - x^7/7 + \dots$. Setting $x = 1$, we have the stated result, since $\tan^{-1} 1 = \pi/4$.

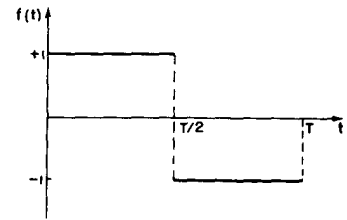


Fig. 50-3. Square-wave function.
 $f(t) = +1$ for $0 < t < T/2$,
 $f(t) = -1$ for $T/2 < t < T$.

We have left only the square terms like $a_5^2 \cos^2 5\omega t$. The integral of any cosine squared or sine squared over one period is equal to $T/2$, so we get

$$\begin{aligned} \int_0^T f^2(t) dt &= Ta_0^2 + \frac{T}{2} (a_1^2 + a_2^2 + \cdots + b_1^2 + b_2^2 + \cdots) \\ &= Ta_0^2 + \frac{T}{2} \sum_{n=1}^{\infty} (a_n^2 + b_n^2). \end{aligned} \quad (50.23)$$

This equation is called the “energy theorem,” and says that the total energy in a wave is just the sum of the energies in all of the Fourier components. For example, applying this theorem to the series (50.19), since $[f(t)]^2 = 1$ we get

$$T = \frac{T}{2} \cdot \left(\frac{4}{\pi}\right)^2 \left(1 + \frac{1}{3^2} + \frac{1}{5^2} + \frac{1}{7^2} \cdots\right),$$

so we learn that the sum of the squares of the reciprocals of the odd integers is $\pi^2/8$. In a similar way, by first obtaining the Fourier series for the function and using the energy theorem, we can prove that $1 + 1/2^4 + 1/3^4 + \cdots$ is $\pi^4/90$, a result we needed in Chapter 45.

50-6 Nonlinear responses

Finally, in the theory of harmonics there is an important phenomenon which should be remarked upon because of its practical importance—that of nonlinear effects. In all the systems that we have been considering so far, we have supposed that everything was linear, that the responses to forces, say the displacements or the accelerations, were always proportional to the forces. Or that the currents in the circuits were proportional to the voltages, and so on. We now wish to consider cases where there is not a strict proportionality. We think, at the moment, of some device in which the response, which we will call x_{out} at the time t , is determined by the input x_{in} at the time t . For example, x_{in} might be the force and x_{out} might be the displacement. Or x_{in} might be the current and x_{out} the voltage. If the device is linear, we would have

$$x_{out}(t) = Kx_{in}(t), \quad (50.24)$$

where K is a constant independent of t and of x_{in} . Suppose, however, that the device is nearly, but not exactly, linear, so that we can write

$$x_{out}(t) = K[x_{in}(t) + \epsilon x_{in}^2(t)], \quad (50.25)$$

where ϵ is small in comparison with unity. Such linear and nonlinear responses are shown in the graphs of Fig. 50-4.

Nonlinear responses have several important practical consequences. We shall discuss some of them now. First we consider what happens if we apply a pure tone at the input. We let $x_{in} = \cos \omega t$. If we plot x_{out} as a function of time we get the solid curve shown in Fig. 50-5. The dashed curve gives, for comparison, the response of a linear system. We see that the output is no longer a cosine function. It is more peaked at the top and flatter at the bottom. We say that the output is *distorted*. We know, however, that such a wave is no longer a pure tone, that it will have harmonics. We can find what the harmonics are. Using $x_{in} = \cos \omega t$ with Eq. (50.25), we have

$$x_{out} = K(\cos \omega t + \epsilon \cos^2 \omega t). \quad (50.26)$$

From the equality $\cos^2 \theta = \frac{1}{2}(1 + \cos 2\theta)$, we have

$$x_{out} = K \left(\cos \omega t + \frac{\epsilon}{2} + \frac{\epsilon}{2} \cos 2\omega t \right). \quad (50.27)$$

The output has not only a component at the fundamental frequency, that was present at the input, but also has some of its second harmonic. There has also

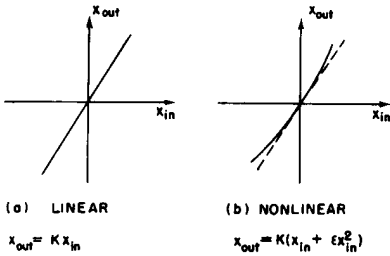


Fig. 50-4. Linear and nonlinear responses.

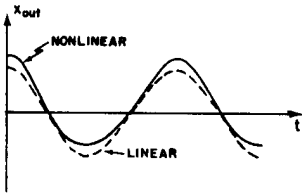


Fig. 50-5. The response of a nonlinear device to the input $\cos \omega t$. A linear response is shown for comparison.

appeared at the output a constant term $K(\epsilon/2)$, which corresponds to the shift of the average value, shown in Fig. 50-5. The process of producing a shift of the average value is called *rectification*.

A nonlinear response will rectify and will produce harmonics of the frequencies at its input. Although the nonlinearity we assumed produced only second harmonics, nonlinearities of higher order—those which have terms like x_{in}^3 and x_{in}^4 , for example—will produce harmonics higher than the second.

Another effect which results from a nonlinear response is *modulation*. If our input function contains two (or more) pure tones, the output will have not only their harmonics, but still other frequency components. Let $x_{in} = A \cos \omega_1 t + B \cos \omega_2 t$, where now ω_1 and ω_2 are *not* intended to be in a harmonic relation. In addition to the linear term (which is K times the input) we shall have a component in the output given by

$$x_{out} = K\epsilon(A \cos \omega_1 t + B \cos \omega_2 t)^2 \quad (50.28)$$

$$= K\epsilon(A^2 \cos^2 \omega_1 t + B^2 \cos^2 \omega_2 t + 2AB \cos \omega_1 t \cos \omega_2 t). \quad (50.29)$$

The first two terms in the parentheses of Eq. (50.29) are just those which gave the constant terms and second harmonic terms we found above. The last term is new.

We can look at this new “cross term” $AB \cos \omega_1 t \cos \omega_2 t$ in two ways. First, if the two frequencies are widely different (for example, if ω_1 is much greater than ω_2) we can consider that the cross term represents a cosine oscillation of varying amplitude. That is, we can think of the factors in this way:

$$AB \cos \omega_1 t \cos \omega_2 t = C(t) \cos \omega_1 t, \quad (50.30)$$

with

$$C(t) = AB \cos \omega_2 t. \quad (50.31)$$

We say that the amplitude of $\cos \omega_1$ is *modulated* with the frequency ω_2 .

Alternatively, we can write the cross term in another way:

$$AB \cos \omega_1 t \cos \omega_2 t = \frac{AB}{2} [\cos (\omega_1 + \omega_2)t + \cos (\omega_1 - \omega_2)t]. \quad (50.32)$$

We would now say that two *new* components have been produced, one at the *sum* frequency $(\omega_1 + \omega_2)$, another at the *difference* frequency $(\omega_1 - \omega_2)$.

We have two different, but equivalent, ways of looking at the same result. In the special case that $\omega_1 \gg \omega_2$, we can relate these two different views by remarking that since $(\omega_1 + \omega_2)$ and $(\omega_1 - \omega_2)$ are near to each other we would expect to observe beats between them. But these beats have just the effect of *modulating* the amplitude of the *average* frequency ω_1 by one-half the difference frequency $2\omega_2$. We see, then, why the two descriptions are equivalent.

In summary, we have found that a nonlinear response produces several effects: rectification, generation of harmonics, and modulation, or the generation of components with sum and difference frequencies.

We should notice that all these effects (Eq. 50.29) are proportional not only to the nonlinearity coefficient ϵ , but also to the product of two amplitudes—either A^2 , B^2 , or AB . We expect these effects to be much more important for *strong* signals than for weak ones.

The effects we have been describing have many practical applications. First, with regard to sound, it is believed that the ear is nonlinear. This is believed to account for the fact that with loud sounds we have the sensation that we *hear* harmonics and also sum and difference frequencies even if the sound waves contain only pure tones.

The components which are used in sound-reproducing equipment—amplifiers, loudspeakers, etc.—always have some nonlinearity. They produce distortions in the sound—they generate harmonics, etc.—which were not present in the original sound. These new components are heard by the ear and are apparently objectionable. It is for this reason that “Hi-Fi” equipment is designed to be as linear as

possible. (Why the nonlinearities of the *ear* are *not* “objectionable” in the same way, or how we even know that the nonlinearity is in the *loudspeaker* rather than in the *ear* is not clear!)

Nonlinearities are quite *necessary*, and are, in fact, intentionally made large in certain parts of radio transmitting and receiving equipment. In an AM transmitter the “voice” signal (with frequencies of some kilocycles per second) is combined with the “carrier” signal (with a frequency of some megacycles per second) in a nonlinear circuit called a *modulator*, to produce the modulated oscillation that is transmitted. In the receiver, the components of the received signal are fed to a nonlinear circuit which combines the sum and difference frequencies of the modulated carrier to generate again the voice signal.

When we discussed the transmission of light, we assumed that the induced oscillations of charges were proportional to the electric field of the light—that the response was linear. That is indeed a very good approximation. It is only within the last few years that light sources have been devised (lasers) which produce an intensity of light strong enough so that nonlinear effects can be observed. It is now possible to generate harmonics of light frequencies. When a strong red light passes through a piece of glass, a little bit of blue light—second harmonic—comes out!

Waves

51-1 Bow waves

Although we have finished our quantitative analyses of waves, this added chapter on the subject is intended to give some appreciation, qualitatively, for various phenomena that are associated with waves, which are too complicated to analyze in detail here. Since we have been dealing with waves for several chapters, more properly the subject might be called "some of the more complex phenomena associated with waves."

The first topic to be discussed concerns the effects that are produced by a source of waves which is moving faster than the wave velocity, or the phase velocity. Let us first consider waves that have a definite velocity, like sound and light. If we have a source of sound which is moving faster than the speed of sound, then something like this happens: Suppose at a given moment a sound wave is generated from the source at point x_1 in Fig. 51-1; then, in the next moment, as the source moves to x_2 , the wave from x_1 expands by a radius r_1 smaller than the distance that the source moves; and, of course, another wave starts from x_2 . When the sound source has moved still farther, to x_3 , and a wave is starting there, the wave from x_2 has now expanded to r_2 , and the one from x_1 has expanded to r_3 . Of course the thing is done continuously, not in steps, and therefore, we have a series of wave circles with a common tangent line which goes through the center of the source. We see that instead of a source generating spherical waves, as it would if it were standing still, it generates a wavefront which forms a cone in three dimensions, or a pair of lines in two dimensions. The angle of the cone is very easy to figure out. In a given amount of time the source moves a distance, say $x_3 - x_1$, proportional to v , the velocity of the source. In the meantime the wavefront has moved out a distance r_3 , proportional to c_w , the speed of the wave. Therefore it is clear that the half-angle of opening has a sine equal to the ratio of the speed of the waves, divided by the speed of the source, and this sine has a solution only if c_w is less than v , or the speed of the object is faster than the speed of the wave:

$$\sin \theta = \frac{c_w}{v}. \quad (51.1)$$

Incidentally, although we implied that it is necessary to have a *source* of sound, it turns out, very interestingly, that once the object is moving faster than the speed of sound, it will *make* sound. That is, it is not necessary that it have a certain tone vibrational character. Any object moving through a medium faster than the speed at which the medium carries waves will generate waves on each side, automatically, just from the motion itself. This is simple in the case of sound, but it also occurs in the case of light. At first one might think nothing can move faster than the speed of light. However, light in glass has a phase velocity less than the speed of light in a vacuum, and it is possible to shoot a charged particle of very high energy through a block of glass such that the particle velocity is close to the speed of light in a vacuum, while the speed of light in the glass may be only $\frac{2}{3}$ the speed of light in the vacuum. A particle moving faster than the speed of light in the medium will produce a conical wave of light with its apex at the source, like the wave wake from a boat (which is from the same effect, as a matter of fact). By measuring the cone angle, we can determine the speed of the particle. This is used technically to determine the speeds of particles as one of the methods of determining their energy in high-energy research. The direction of the light is all that needs to be measured.

51-1 Bow waves

51-2 Shock waves

51-3 Waves in solids

51-4 Surface waves

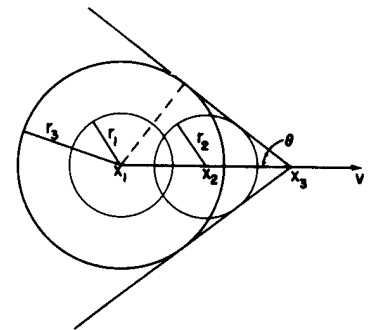


Fig. 51-1. The shock wave front lies on a cone with apex at the source and half-angle $\theta = \sin^{-1} v/c_w$.

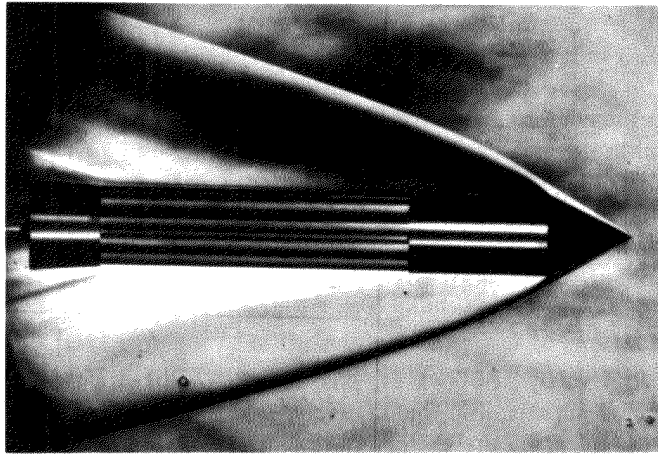


Fig. 51-2. A shock wave induced in a gas by a projectile moving faster than sound.

This light is sometimes called Čerenkov radiation, because it was first observed by Čerenkov. How intense this light should be was analyzed theoretically by Frank and Tamm. The 1958 Nobel Prize for physics was awarded jointly to all three for this work.

The corresponding circumstances in the case of sound are illustrated in Fig. 51-2, which is a photograph of an object moving through a gas at a speed greater than the speed of sound. The changes in pressure produce a change in refractive index, and with a suitable optical system the edges of the waves can be made visible. We see that the object moving faster than the speed of sound does, indeed, produce a conical wave. But closer inspection reveals that the surface is actually curved. It is straight asymptotically, but it is curved near the apex, and we have now to discuss how that can be, which brings us to the second topic of this chapter.

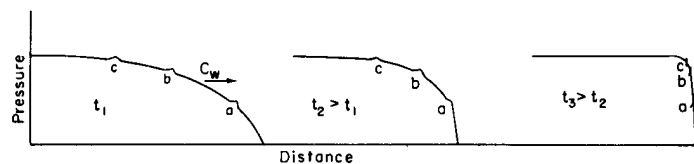


Fig. 51-3. Wavefront "snapshots" at successive instants in time.

51-2 Shock waves

Wave speed often depends on the amplitude, and in the case of sound the speed depends upon the amplitude in the following way. An object moving through the air has to move the air out of the way, so the disturbance produced in this case is some kind of a pressure step, with the pressure higher behind the wavefront than in the undisturbed region not yet reached by the wave (running along at the normal speed, say). But the air that is left behind, after the wavefront passes, has been compressed adiabatically, and therefore the temperature is increased. Now the speed of sound increases with the temperature, so the speed in the region behind the jump is faster than in the air in front. That means that any other disturbance that is made behind this step, say by a continuous pushing of the body, or any other disturbance, will ride faster than the front, the speed increasing with higher pressure. Figure 51-3 illustrates the situation, with some little bumps of pressure added to the pressure contour to aid visualization. We see that the higher pressure regions at the rear overtake the front as time goes on, until ultimately the compressional wave develops a sharp front. If the strength is very high, "ultimately" means right away; if it is rather weak, it takes a long time; it may be, in fact, that the sound is spreading and dying out before it has time to do this.

The sounds we make in talking are extremely weak relative to the atmospheric pressure—only 1 part in a million or so. But for pressure changes of the order of 1 atmosphere, the wave velocity increases by about twenty percent, and the wavefront sharpens up at a correspondingly high rate. In nature nothing happens *infinitely* rapidly, presumably, and what we call a "sharp" front has, actually, a very slight thickness; it is not infinitely steep. The distances over which it is varying are of the

order of one mean free path, in which the theory of the wave equation begins to fail because we did not consider the structure of the gas.

Now, referring again to Fig. 51-2, we see that the curvature can be understood if we appreciate that the pressures near the apex are higher than they are farther back, and so the angle θ is greater. That is, the curve is the result of the fact that the speed depends upon the strength of the wave. Therefore the wave from an atomic bomb explosion travels much faster than the speed of sound for a while, until it gets so far out that it is weakened to such an extent from spreading that the pressure bump is small compared with atmospheric pressure. The speed of the bump then approaches the speed of sound in the gas into which it is going. (Incidentally, it always turns out that the speed of the shock is higher than the speed of sound in the gas ahead, but is lower than the speed of sound in the gas behind. That is, impulses from the back will arrive at the front, but the front rides into the medium in which it is going faster than the normal speed of signals. So one cannot tell, acoustically, that the shock is coming until it is too late. The light from the bomb arrives first, but one cannot tell that the shock is coming until it arrives, because there is no sound signal coming ahead of it.)

This is a very interesting phenomenon, this piling up of waves, and the main point on which it depends is that after a wave is present, the speed of the resulting wave should be higher. Another example of the same phenomenon is the following. Consider water flowing in a long channel with finite width and finite depth. If a piston, or a wall across the channel, is moved along the channel fast enough, water piles up, like snow before a snow plow. Now suppose the situation is as shown in Fig. 51-4, with a sudden step in water height somewhere in the channel. It can be demonstrated that long waves in a channel travel faster in deeper water than they do in shallow water. Therefore any new bumps or irregularities in energy supplied by the piston run off forward and pile up at the front. Again, ultimately what we have is just water with a sharp front, theoretically. However, as Fig. 51-4 shows, there are complications. Pictured is a wave coming up a channel; the piston is at the far right end of the channel. At first it might have appeared like a well-behaved wave, as one might expect, but farther along the channel, it has become sharper and sharper until the events pictured occurred. There is a terrible churning at the surface, as the pieces of water fall down, but it is essentially a very sharp rise with no disturbance of the water ahead.

Actually water is much more complicated than sound. However, just to illustrate a point, we will try to analyze the speed of such a so-called *bore*, in a channel. The point here is not that this is of any basic importance for our purposes—it is not a great generalization—it is only to illustrate that the laws of mechanics that we already know are capable of explaining the phenomenon.

Imagine, for a moment, that the water does look something like Fig. 51-5(a), that water at the higher height h_2 is moving with a velocity v , and that the front is moving with velocity u into undisturbed water which is at height h_1 . We would like to determine the speed at which the front moves. In a time Δt a vertical plane initially at x_1 moves a distance $v \Delta t$ to x_2 , while the front of the wave has moved $u \Delta t$.

Now we apply the equations of conservation of matter and momentum. First, the former: Per unit channel width, we see that the amount $h_2 v \Delta t$ of matter that has moved past x_1 (shown shaded) is compensated by the other shaded region, which amounts to $(h_2 - h_1)u \Delta t$. So, dividing by Δt , $vh_2 = u(h_2 - h_1)$. That does not yet give us enough, because although we have h_2 and h_1 , we do not know either u or v ; we are trying to get both of them.

Now the next step is to use conservation of momentum. We have not discussed the problems of water pressure, or anything in hydrodynamics, but it is clear anyway that the pressure of water at a given depth is just enough to hold up the column of water above it. Therefore the pressure of water is equal to ρ , the density of water, times g , times the depth below the surface. Since the pressure increases linearly with depth, the average pressure over the plane at x_1 , say, is $\frac{1}{2}\rho gh_2$, which is also the average force per unit width and per unit height pushing the plane toward x_2 . So we multiply by another h_2 to get the total force which is acting on the water



Figure 51-4

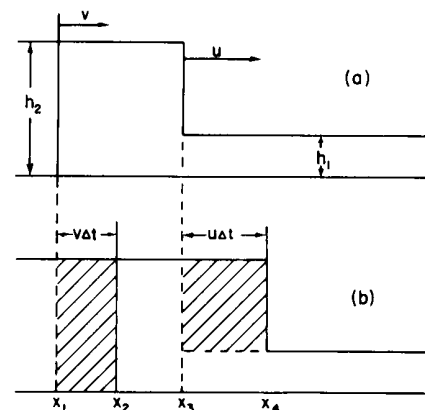


Fig. 51-5. Two cross sections of a bore in a channel, with (b) at an interval Δt later than (a).

pushing from the left. On the other hand, there is pressure in the water on the right also, exerting an opposite force on the region in question, which is, by the same kind of analysis, $\frac{1}{2}\rho gh_1^2$. Now we must balance the forces against the rate of change of the momentum. Thus we have to figure out how much more momentum there is in situation (b) in Fig. 51-5 than there was in (a). We see that the additional mass that has acquired the speed v is just $\rho h_2 u \Delta t - \rho h_2 v \Delta t$ (per unit width), and multiplying this by v gives the additional momentum to be equated to the impulse $F \Delta t$:

$$(\rho h_2 u \Delta t - \rho h_2 v \Delta t)v = (\frac{1}{2}\rho gh_2^2 - \frac{1}{2}\rho gh_1^2) \Delta t.$$

If we eliminate v from this equation by substituting $vh_2 = u(h_2 - h_1)$, already found, and simplify, we get finally that $u^2 = gh_2(h_1 + h_2)/2h_1$.

If the height difference is very small, so that h_1 and h_2 are nearly equal, this says that the velocity $= \sqrt{gh}$. As we will see later, that is only true provided the wavelength of the wave is longer than the depth of the channel.

We could also do the analogous thing for sound waves—including the conservation of internal energy, not the conservation of entropy, because the shock is irreversible. In fact, if one checks the conservation of energy in the bore problem, one finds that energy is not conserved. If the height difference is small, it is almost perfectly conserved, but as soon as the height difference becomes very appreciable, there is a net loss of energy. This is manifested as the falling water and the churning shown in Fig. 51-4.

In shock waves there is a corresponding apparent loss of energy, from the point of view of adiabatic reactions. The energy in the sound wave, behind the shock, goes into heating of the gas after shock passes, corresponding to churning of the water in the bore. In working it out, three equations for the sound case turn out to be necessary for solution, and the temperature behind the shock is not the same as the temperature in front, as we have seen.

If we try to make a bore that is upside down ($h_2 < h_1$), then we find that the energy loss per second is negative. Since energy is not available from anywhere, that bore cannot then maintain itself; it is unstable. If we were to start a wave of that sort, it would flatten out, because the speed dependence on height that resulted in sharpening in the case we discussed would now have the opposite effect.

51-3 Waves in solids

The next kind of waves to be discussed are the more complicated waves in solids. We have already discussed sound waves in gas and in liquid, and there is a direct analog to a sound wave in a solid. If a sudden push is applied to a solid, it is compressed. It resists the compression, and a wave analogous to sound is started. However there is another kind of wave that is possible in a solid, and which is not possible in a fluid. If a solid is distorted by pushing it sideways (called *shearing*), then it tries to pull itself back. That is by definition what distinguishes a solid from a liquid: if we distort a liquid (internally), hold it a minute so that it calms down, and then let go, it will stay that way, but if we take a solid and push it, like shearing a piece of "Jello," and let it go, it flies back and starts a *shear* wave, travelling in the same way the compressions travel. In all cases, the shear wave speed is less than the speed of longitudinal waves. The shear waves are somewhat more analogous, so far as their polarizations are concerned, to light waves. Sound has no polarization, it is just a pressure wave. Light has a characteristic orientation perpendicular to its direction of travel.

In a solid, the waves are of both kinds. First, there is a compression wave, analogous to sound, that runs at one speed. If the solid is not crystalline, then a shear wave polarized in any direction will propagate at a characteristic speed. (Of course all solids are crystalline, but if we use a block made up of microcrystals of all orientations, the crystal anisotropies average out.)

Another interesting question concerning sound waves is the following: What happens if the wavelength in a solid gets shorter, and shorter, and shorter? How short can it get? It is interesting that it cannot get any shorter than the space

between the atoms, because if there is supposed to be a wave in which one point goes up and the next down, etc., the shortest possible wavelength is clearly the atom spacing. In terms of the modes of oscillation, we say that there are longitudinal modes, and transverse modes, long wave modes, short wave modes. As we consider wavelengths comparable to the spacing between the atoms, then the speeds are no longer constant; there is a dispersion effect where the velocity is not independent of the wave number. But, ultimately, the highest mode of transverse waves would be that in which every atom is doing the opposite of neighboring atoms.

Now from the point of view of atoms, the situation is like the two pendulums that we were talking about, for which there are two modes, one in which they both go together, and the other in which they go apart. It is possible to analyze the solid waves another way, in terms of a system of coupled harmonic oscillators, like an enormous number of pendulums, with the highest mode such that they oscillate oppositely, and lower modes with different relationships of the timing.

The shortest wavelengths are so short that they are not usually available technically. However they are of great interest because, in the theory of thermodynamics of a solid, the heat properties of a solid, for example specific heats, can be analyzed in terms of the properties of the short sound waves. Going to the extreme of sound waves of ever shorter wavelength, one necessarily comes to the individual motions of the atoms; the two things are the same ultimately.

A very interesting example of sound waves in a solid, both longitudinal and transverse, are the waves that are in the solid earth. Who makes the noises we do not know, but inside the earth, from time to time, there are earthquakes—some rock slides past some other rock. That is like a little noise. So waves like sound waves start out from such a source very much longer in wavelength than one usually considers in sound waves, but still they are sound waves, and they travel around in the earth. The earth is not homogeneous, however, and the properties of pressure, density, compressibility, and so on, change with depth, and therefore the speed varies with depth. Then the waves do not travel in straight lines—there is a kind of index of refraction and they go in curves. The longitudinal waves and the transverse waves have different speeds, so there are different solutions for the different speeds. Therefore if we place a seismograph at some location and watch the way the thing jiggles after there has been an earthquake somewhere else, then we do not just get an irregular jiggling. We might get a jiggling, and a quieting down, and then another jiggling—what happens depends upon the location. If it were close enough, we would first receive longitudinal waves from the disturbance, and then, a few moments later, transverse waves, because they travel more slowly. By measuring the time difference between the two, we can tell how far away the earthquake is, if we know enough about the speeds and composition of the interior regions involved.

An example of the behavior pattern of waves in the earth is shown in Fig. 51-6. The two kinds of waves are represented by different symbols. If there were an earthquake at the place marked "source," the transverse waves and longitudinal waves would arrive at different times at the station by the most direct routes, and there would also be reflections at discontinuities, resulting in other paths and times. It turns out that there is a core in the earth which does not carry transverse waves. If the station is opposite the source, transverse waves still arrive, but the timing is not right. What happens is that the transverse wave comes to the core, and whenever the transverse waves come to a surface which is oblique, between two materials, two new waves are generated, one transverse and one longitudinal. But inside the core of the earth, a transverse wave is not propagated (or at least, there is no evidence for it, only for a longitudinal wave); it comes out again in both forms and comes to the station.

It is from the behavior of these earthquake waves that it has been determined that transverse waves cannot be propagated within the inner circle. This means that the center of the earth is liquid in the sense that it cannot propagate transverse waves. The only way we know what is inside the earth is by studying earthquakes. So, by using a large number of observations of many earthquakes at different

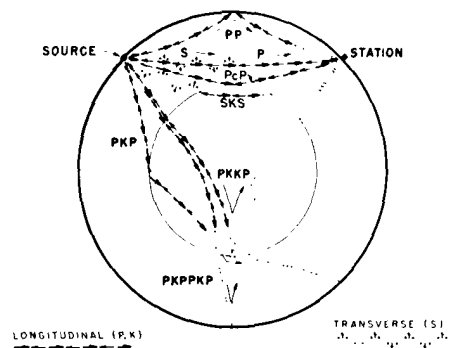


Fig. 51-6. Schematic of the earth, showing paths of longitudinal and transverse sound waves.

stations, the details have been worked out—the speed, the curves, etc. are all known. We know what the speeds of various kinds of waves are at every depth. Knowing that, therefore, it is possible to figure out what the normal modes of the earth are, because we know the speed of propagation of sound waves—in other words, the elastic properties of both kinds of waves at every depth. Suppose the earth were distorted into an ellipsoid and let go. It is just a matter of superposing waves travelling around in the ellipsoid to determine the period and shapes in a free mode. We have figured out that if there is a disturbance, there are a lot of modes, from the lowest, which is ellipsoidal, to higher modes with more structure.

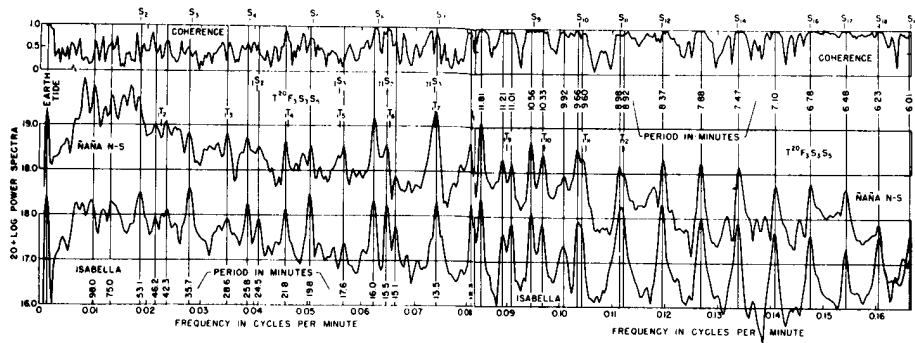


Fig. 51-7. Power versus frequency as detected at seismographs in Naña, Peru, and Isabella, California. The coherence is a measure of the coupling between the stations. [From Benioff, Press and Smith, *J. Geoph. Research* **66**, 605 (1961)].

The Chilean earthquake of May 1960 made a loud enough “noise” that the signals went around the earth many times, and new seismographs of great delicacy were made just in time to determine the frequencies of the fundamental modes of the earth and to compare them with the values that were calculated from the theory of sound with the known velocities, as measured from the independent earthquakes. The result of this experiment is illustrated in Fig. 51-7, which is a plot of the strength of the signal versus the frequency of its oscillation (a *Fourier analysis*). Note that at certain particular frequencies there is much more being received than at other frequencies; there are very definite maxima. These are the natural frequencies of the earth, because these are the main frequencies at which the earth can oscillate. In other words, if the entire motion of the earth is made up of many different modes, we would expect to obtain, for each station, irregular bumpings which indicate a superposition of many frequencies. If we analyze this in terms of frequencies, we should be able to find the characteristic frequencies of the earth. The vertical dark lines in the figure are the calculated frequencies, and we find a remarkable agreement, an agreement due to the fact that the theory of sound is right for the inside of the earth.

A very curious point is revealed in Fig. 51-8, which shows a very careful measurement, with better resolution of the lowest mode, the ellipsoidal mode of the earth. Note that it is not a single maximum, but a double one, 54.7 minutes and 53.1 minutes—slightly different. The reason for the two different frequencies was not known at the time that it was measured, although it may have been found in the meantime. There are at least two possible explanations: One would be that there may be asymmetry in the earth’s distribution, which would result in two similar modes. Another possibility, which is even more interesting, is this: Imagine the waves going around the earth in two directions from the source. The speeds will not be equal because of effects of the rotation of the earth in the equations of motion, which have not been taken into account in making the analysis. Motion in a rotating system is modified by Coriolis forces, and these may cause the observed splitting.

Regarding the method by which these quakes have been analyzed, what is obtained on the seismograph is not a curve of amplitude as a function of frequency, but displacement as a function of time, always a very irregular tracing. To find the amount of all the different sine waves for all different frequencies, we know that the trick is to multiply the data by a sine wave of a given frequency and integrate, i.e., average it, and in the average all other frequencies disappear. The figures were thus plots of the integrals found when the data were multiplied by sine waves of different cycles per minute, and integrated.

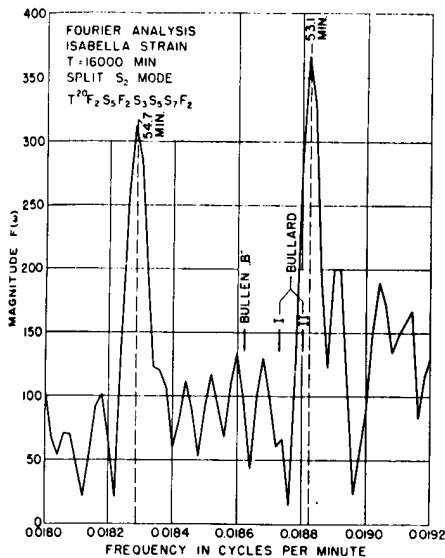
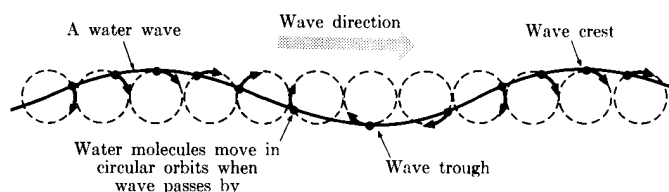


Fig. 51-8. High-resolution analysis of one of the seismograph records, showing spectral doublet.

51-4 Surface waves

Now, the next waves of interest, that are easily seen by everyone and which are usually used as an example of waves in elementary courses, are water waves. As we shall soon see, they are the worst possible example, because they are in no respects like sound and light; they have all the complications that waves can have. Let us start with long water waves in deep water. If the ocean is considered infinitely deep and a disturbance is made on the surface, waves are generated. All kinds of irregular motions occur, but the sinusoidal type motion, with a very small disturbance, might look like the common smooth ocean waves coming in toward the shore. Now with such a wave, the water, of course, on the average, is standing still, but the wave moves. What is the motion, is it transverse or longitudinal? It must be neither; it is not transverse, nor is it longitudinal. Although the water at a given place is alternately trough or hill, it cannot simply be moving up and down, by the conservation of water. That is, if it goes down, where is the water going to go? The water is essentially incompressible. The speed of compression of waves—that is, sound in the water—is much, much higher, and we are not considering that now. Since water is incompressible on this scale, as a hill comes down the water must move away from the region. What actually happens is that particles of water near the surface move approximately in circles. When smooth swells are coming, a person floating in a tire can look at a nearby object and see it going in a circle. So it is a mixture of longitudinal and transverse, to add to the confusion. At greater depths in the water the motions are smaller circles until, reasonably far down, there is nothing left of the motion (Fig. 51-9).

Fig. 51-9. Deep-water waves are formed from particles moving in circles. Note the systematic phase shifts from circle to circle. How would a floating object move?



To find the velocity of such waves is an interesting problem: it must be some combination of the density of the water, the acceleration of gravity, which is the restoring force that makes the waves, and possibly of the wavelength and of the depth. If we take the case where the depth goes to infinity, it will no longer depend on the depth. Whatever formula we are going to get for the velocity of the phases of the waves must combine the various factors to make the proper dimensions, and if we try this in various ways, we find only one way to combine the density, g , and λ in order to make a velocity, namely, $\sqrt{g\lambda}$, which does not include the density at all. Actually, this formula for the phase velocity is not exactly right, but a complete analysis of the dynamics, which we will not go into, shows that the factors are as we have them, except for $\sqrt{2\pi}$:

$$v_{\text{phase}} = \sqrt{g\lambda/2\pi} \text{ (for gravity waves).}$$

It is interesting that the long waves go faster than the short waves. Thus if a boat makes waves far out, because there is some sports-car driver in a motorboat travelling by, then after a while the waves come to shore with slow sloshings at first and then more and more rapid sloshings, because the first waves that come are long. The waves get shorter and shorter as the time goes on, because the velocities go as the square root of the wavelength.

One may object, "That is not right, we must look at the *group* velocity in order to figure it out!" Of course that is true. The formula for the phase velocity does not tell us what is going to arrive first; what tells us is the group velocity. So we have to work out the group velocity, and it is left as a problem to show it to be one-half of the phase velocity, assuming that the velocity goes as the square root of the wavelength, which is all that is needed. The group velocity also goes as the square root of the wavelength. How can the group velocity go half as fast as the phase? If one looks at the bunch of waves that are made by a boat travelling

along, following a particular crest, he finds that it moves forward in the group and gradually gets weaker and dies out in the front, and mystically and mysteriously a weak one in the back works its way forward and gets stronger. In short, the waves are moving through the group while the group is only moving at half the speed that the waves are moving.

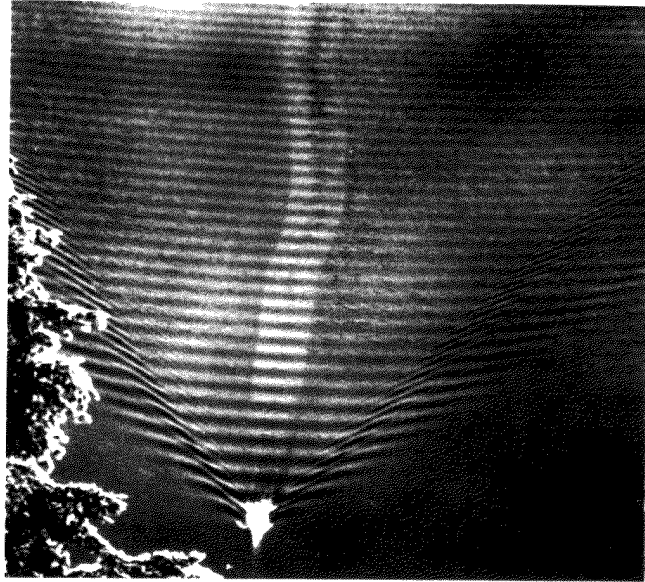


Fig. 51-10. The wake of a boat.

Because the group velocities and phase velocities are not equal, then the waves that are produced by an object moving through are no longer simply a cone, but it is much more interesting. We can see that in Fig. 51-10, which shows the waves produced by an object moving through the water. Note that it is quite different than what we would have for sound, in which the velocity is independent of wavelength, where we would have wavefronts only along the cone, travelling outward. Instead of that, we have waves in the back with fronts moving parallel to the motion of the boat, and then we have little waves on the sides at other angles. This entire pattern of waves can, with ingenuity, be analyzed by knowing only this: that the phase velocity is proportional to the square root of the wavelength. The trick is that the pattern of waves is stationary relative to the (constant-velocity) boat; any other pattern would get lost from the boat.

The water waves that we have been considering so far were long waves in which the force of restoration is due to gravitation. But when waves get very short in the water, the main restoring force is capillary attraction, i.e., the energy of the surface, the surface tension. For surface tension waves, it turns out that the phase velocity is

$$v_{\text{phase}} = \sqrt{2\pi T / \lambda \rho} \text{ (for ripples),}$$

where T is the surface tension and ρ the density. It is the exact opposite: the phase velocity is *higher*, the shorter the wavelength, when the wavelength gets very small. When we have both gravity and capillary action, as we always do, we get the combination of these two together:

$$v_{\text{phase}} = \sqrt{Tk/\rho + g/k},$$

where $k = 2\pi/\lambda$ is the wave number. So the velocity of the waves of water is really quite complicated. The phase velocity as a function of the wavelength is shown in Fig. 51-11; for very short waves it is fast, for very long waves it is fast, and there is a minimum speed at which the waves can go. The group velocity can be calculated from the formula: it goes to $\frac{3}{2}$ the phase velocity for ripples and $\frac{1}{2}$ the phase velocity for gravity waves. To the left of the minimum the group velocity is higher than the phase velocity; to the right, the group velocity is less than the

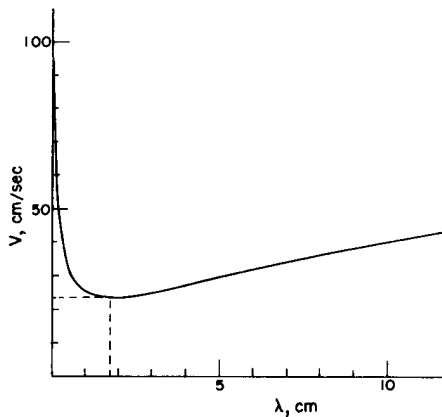


Fig. 51-11. Phase velocity vs. wavelength for water.

phase velocity. There are a number of interesting phenomena associated with these facts. In the first place, since the group velocity is increasing so rapidly as the wavelength goes down, if we make a disturbance there will be a slowest end of the disturbance going at the minimum speed with the corresponding wavelength, and then in front, going at higher speed, will be a short wave and a very long wave. It is very hard to see the long ones, but it is easy to see the short ones in a water tank.

So we see that the ripples often used to illustrate simple waves are quite interesting and complicated; they do not have a sharp wavefront at all, as is the case for simple waves like sound and light. The main wave has little ripples which run out ahead. A sharp disturbance in the water does not produce a sharp wave because of the dispersion. First come the very fine waves. Incidentally, if an object moves through the water at a certain speed, a rather complicated pattern results, because all the different waves are going at different speeds. One can demonstrate this with a tray of water and see that the fastest ones are the fine capillary waves. There are slowest waves, of a certain kind, which go behind. By inclining the bottom, one sees that where the depth is lower, the speed is lower. If a wave comes in at an angle to the line of maximum slope, it bends and tends to follow that line. In this way one can show various things, and we conclude that waves are more complicated in water than in air.

The speed of long waves in water with circulatory motions is slower when the depth is less, faster in deep water. Thus as water comes toward a beach where the depth lessens, the waves go slower. But where the water is deeper, the waves are faster, so we get the effects of shock waves. This time, since the wave is not so simple, the shocks are much more contorted, and the wave over-curves itself, in the familiar way shown in Fig. 51-12. This is what happens when waves come into the shore, and the real complexities in nature are well revealed in such a circumstance. No one has yet been able to figure out what shape the wave should take as it breaks. It is easy enough when the waves are small, but when one gets large and breaks, then it is much more complicated.



Fig. 51-12. A water wave.

An interesting feature about capillary waves can be seen in the disturbances made by an object moving through the water. From the point of view of the object itself, the water is flowing past, and the waves which ultimately sit around it are always the waves which have just the right speed to stay still with the object in the water. Similarly, around an object in a stream, with the stream flowing by, the pattern of waves is stationary, and at just the right wavelengths to go at the same speed as the water going by. But if the group velocity is less than the phase velocity, then the disturbances propagate out backwards in the stream, because the group velocity is not quite enough to keep up with the stream. If the group velocity is faster than the velocity of the phase, the pattern of waves will appear in front of the object. If one looks closely at objects in a stream, one can see that there are little ripples in front and long “slurps” in the back.

Another interesting feature of this sort can be observed in pouring liquids. If milk is poured fast enough out of a bottle, for instance, a large number of lines can be seen crossing both ways in the outgoing stream. They are waves starting from the disturbance at the edges and running out, much like the waves about an object in a stream. There are effects from both sides which produce the crossed pattern.

We have investigated some of the interesting properties of waves and the various complications of dependence of phase velocity on wavelength, the speed of the waves on depth, and so forth, that produce the really complex, and therefore interesting, phenomena of nature.

Symmetry in Physical Laws

52-1 Symmetry operations

The subject of this chapter is what we may call *symmetry in physical laws*. We have already discussed certain features of symmetry in physical laws in connection with vector analysis (Chapter 11), the theory of relativity (Chapter 16), and rotation (Chapter 20).

Why should we be concerned with symmetry? In the first place, symmetry is fascinating to the human mind, and everyone likes objects or patterns that are in some way symmetrical. It is an interesting fact that nature often exhibits certain kinds of symmetry in the objects we find in the world around us. Perhaps the most symmetrical object imaginable is a sphere, and nature is full of spheres—stars, planets, water droplets in clouds. The crystals found in rocks exhibit many different kinds of symmetry, the study of which tells us some important things about the structure of solids. Even the animal and vegetable worlds show some degree of symmetry, although the symmetry of a flower or of a bee is not as perfect or as fundamental as is that of a crystal.

But our main concern here is not with the fact that the *objects* of nature are often symmetrical. Rather, we wish to examine some of the even more remarkable symmetries of the universe—the symmetries that exist in the *basic laws themselves* which govern the operation of the physical world.

First, what *is* symmetry? How can a physical *law* be “symmetrical”? The problem of defining symmetry is an interesting one and we have already noted that Weyl gave a good definition, the substance of which is that a thing is symmetrical if there is something we can do to it so that after we have done it, it looks the same as it did before. For example, a symmetrical vase is of such a kind that if we reflect or turn it, it will look the same as it did before. The question we wish to consider here is what we can do to physical phenomena, or to a physical situation in an experiment, and yet leave the result the same. A list of the known operations under which various physical phenomena remain invariant is shown in Table 52-1.

52-2 Symmetry in space and time

The first thing we might try to do, for example, is to *translate* the phenomenon in space. If we do an experiment in a certain region, and then build another apparatus at another place in space (or move the original one over) then, whatever went on in one apparatus, in a certain order in time, will occur in the same way if we have arranged the same condition, with all due attention to the restrictions that we mentioned before: that all of those features of the environment which make it not behave the same way have also been moved over—we talked about how to define how much we should include in those circumstances, and we shall not go into those details again.

In the same way, we also believe today that *displacement in time* will have no effect on physical laws. (That is, *as far as we know today*—all of these things are as far as we know today!) That means that if we build a certain apparatus and start it at a certain time, say on Thursday at 10:00 a.m., and then build the same apparatus and start it, say, three days later in the same condition, the two apparatuses will go through the same motions in exactly the same way as a function of time no matter what the starting time, provided again, of course, that the relevant features of the environment are also modified appropriately in *time*. That symmetry means,

52-1 Symmetry operations

52-2 Symmetry in space and time

52-3 Symmetry and conservation laws

52-4 Mirror reflections

52-5 Polar and axial vectors

52-6 Which hand is right?

52-7 Parity is not conserved!

52-8 Antimatter

52-9 Broken symmetries

Table 52-1

Symmetry Operations

Translation in space
Translation in time
Rotation through a fixed angle
Uniform velocity in a straight line (Lorentz transformation)
Reversal of time
Reflection of space
Interchange of identical atoms or identical particles
Quantum-mechanical phase
Matter-antimatter (charge conjugation)

of course, that if one bought General Motors stock three months ago, the same thing would happen to it if he bought it now!

We have to watch out for geographical differences too, for there are, of course, variations in the characteristics of the earth's surface. So, for example, if we measure the magnetic field in a certain region and move the apparatus to some other region, it may not work in precisely the same way because the magnetic field is different, but we say that is because the magnetic field is associated with the earth. We can imagine that if we move the whole earth and the equipment, it would make no difference in the operation of the apparatus.

Another thing that we discussed in considerable detail was rotation in space: if we turn an apparatus at an angle it works just as well, provided we turn everything else that is relevant along with it. In fact, we discussed the problem of symmetry under rotation in space in some detail in Chapter 11, and we invented a mathematical system called *vector analysis* to handle it as neatly as possible.

On a more advanced level we had another symmetry—the symmetry under uniform velocity in a straight line. That is to say—a rather remarkable effect—that if we have a piece of apparatus working a certain way and then take the same apparatus and put it in a car, and move the whole car, plus all the relevant surroundings, at a uniform velocity in a straight line, then so far as the phenomena inside the car are concerned there is no difference: all the laws of physics appear the same. We even know how to express this more technically, and that is that the mathematical equations of the physical laws must be unchanged under a *Lorentz transformation*. As a matter of fact, it was a study of the relativity problem that concentrated physicists' attention most sharply on symmetry in physical laws.

Now the above-mentioned symmetries have all been of a geometrical nature, time and space being more or less the same, but there are other symmetries of a different kind. For example, there is a symmetry which describes the fact that we can replace one atom by another of the same kind; to put it differently, there *are* atoms of the same kind. It is possible to find groups of atoms such that if we change a pair around, it makes no difference—the atoms are identical. Whatever one atom of oxygen of a certain type will do, another atom of oxygen of that type will do. One may say, "That is ridiculous, that is the *definition* of equal types!" That may be merely the definition, but then we still do not know whether there *are* any "atoms of the same type"; the *fact* is that there are many, many atoms of the same type. Thus it does mean something to say that it makes no difference if we replace one atom by another of the same type. The so-called elementary particles of which the atoms are made are also identical particles in the above sense—all electrons are the same; all protons are the same; all positive pions are the same; and so on.

After such a long list of things that can be done without changing the phenomena, one might think we could do practically anything; so let us give some examples to the contrary, just to see the difference. Suppose that we ask: "Are the physical laws symmetrical under a change of scale?" Suppose we build a certain piece of apparatus, and then build another apparatus five times bigger in every part, will it work exactly the same way? The answer is, in this case, *no*! The wavelength of light emitted, for example, by the atoms inside one box of sodium atoms and the wavelength of light emitted by a gas of sodium atoms five times in volume is not five times longer, but is, in fact, exactly the same as the other. So the ratio of the wavelength to the size of the emitter will change.

Another example: we see in the newspaper, every once in a while pictures of a great cathedral made with little matchsticks—a tremendous work of art by some retired fellow who keeps gluing matchsticks together. It is much more elaborate and wonderful than any real cathedral. If we imagine that this wooden cathedral were actually built on the scale of a real cathedral, we see where the trouble is; it would not last—the whole thing would collapse because of the fact that scaled-up matchsticks are just not strong enough. "Yes," one might say, "but we also know that when there is an influence from the outside, it also must be changed in proportion!" We are talking about the ability of the object to withstand gravitation. So what we should do is first to take the model cathedral of real matchsticks and

the real earth, and then we know it is stable. Then we should take the larger cathedral and take a bigger earth. But then it is even worse, because the gravitation is increased still more!

Today, of course, we understand the fact that phenomena depend on the scale on the grounds that matter is atomic in nature, and certainly if we built an apparatus that was so small there were only five atoms in it, it would clearly be something we could not scale up and down arbitrarily. The scale of an individual atom is not at all arbitrary—it is quite definite.

The fact that the laws of physics are not unchanged under a change of scale was discovered by Galileo. He realized that the strengths of materials were not in exactly the right proportion to their sizes, and he illustrated this property that we were just discussing, about the cathedral of matchsticks, by drawing two bones, the bone of one dog, in the right proportion for holding up his weight, and the imaginary bone of a “super dog” that would be, say, ten or a hundred times bigger—that bone was a big, solid thing with quite different proportions. We do not know whether he ever carried the argument quite to the conclusion that the laws of nature must have a definite scale, but he was so impressed with this discovery that he considered it to be as important as the discovery of the laws of motion, because he published them both in the same volume, called “On Two New Sciences.”

Another example in which the laws are not symmetrical, that we know quite well, is this: a system in rotation at a uniform angular velocity does not give the same apparent laws as one that is not rotating. If we make an experiment and then put everything in a space ship and have the space ship spinning in empty space, all alone at a constant angular velocity, the apparatus will not work the same way because, as we know, things inside the equipment will be thrown to the outside, and so on, by the centrifugal or coriolis forces, etc. In fact, we can tell that the earth is rotating by using a so-called Foucault pendulum, without looking outside.

Next we mention a very interesting symmetry which is obviously false, i.e., *reversibility in time*. The physical laws apparently cannot be reversible in time, because, as we know, all obvious phenomena are irreversible on a large scale: “The moving finger writes, and having writ, moves on.” So far as we can tell, this irreversibility is due to the very large number of particles involved, and if we could see the individual molecules, we would not be able to discern whether the machinery was working forward or backwards. To make it more precise: we build a small apparatus in which we know what all the atoms are doing, in which we can watch them jiggling. Now we build another apparatus like it, but which starts its motion in the final condition of the other one, with all the velocities precisely reversed. *It will then go through the same motions, but exactly in reverse.* Putting it another way: if we take a motion picture, with sufficient detail, of all the inner works of a piece of material and shine it on a screen and run it backwards, no physicist will be able to say, “That is against the laws of physics, that is doing something wrong!” If we do not see all the details, of course, the situation will be perfectly clear. If we see the egg splattering on the sidewalk and the shell cracking open, and so on, then we will surely say, “That is irreversible, because if we run the moving picture backwards the egg will all collect together and the shell will go back together, and that is obviously ridiculous!” But if we look at the individual atoms themselves, the laws look completely reversible. This is, of course, a much harder discovery to have made, but apparently it is true that the fundamental physical laws, on a microscopic and fundamental level, are completely reversible in time!

52-3 Symmetry and conservation laws

The symmetries of the physical laws are very interesting at this level, but they turn out, in the end, to be even more interesting and exciting when we come to quantum mechanics. For a reason which we cannot make clear at the level of the present discussion—a fact that most physicists still find somewhat staggering, a most profound and beautiful thing, is that, in quantum mechanics, *for each of the rules of symmetry there is a corresponding conservation law*; there is a definite

connection between the laws of conservation and the symmetries of physical laws. We can only state this at present, without any attempt at explanation.

The fact, for example, that the laws are symmetrical for translation in space when we add the principles of quantum mechanics, turns out to mean that *momentum is conserved*.

That the laws are symmetrical under translation in time means, in quantum mechanics, that *energy is conserved*.

Invariance under rotation through a fixed angle in space corresponds to the *conservation of angular momentum*. These connections are very interesting and beautiful things, among the most beautiful and profound things in physics.

Incidentally, there are a number of symmetries which appear in quantum mechanics which have no classical analog, which have no method of description in classical physics. One of these is as follows: If ψ is the amplitude for some process or other, we know that the absolute square of ψ is the probability that the process will occur. Now if someone else were to make his calculations, not with this ψ , but with a ψ' which differs merely by a change in phase (let Δ be some constant, and multiply $e^{i\Delta}$ times the old ψ), the absolute square of ψ' , which is the probability of the event, is then equal to the absolute square of ψ :

$$\psi' = \psi e^{i\Delta}; \quad |\psi'|^2 = |\psi|^2. \quad (52.1)$$

Therefore the physical laws are unchanged if the phase of the wave function is shifted by an arbitrary constant. That is another symmetry. Physical laws must be of such a nature that a shift in the quantum-mechanical phase makes no difference. As we have just mentioned, in quantum mechanics there is a conservation law for every symmetry. The conservation law which is connected with the quantum-mechanical phase seems to be the *conservation of electrical charge*. This is altogether a very interesting business!

52-4 Mirror reflections

Now the next question, which is going to concern us for most of the rest of this chapter, is the question of symmetry under *reflection in space*. The problem is this: Are the physical laws symmetrical under reflection? We may put it this way: Suppose we build a piece of equipment, let us say a clock, with lots of wheels and hands and numbers; it ticks, it works, and it has things wound up inside. We look at the clock in the mirror. How it *looks* in the mirror is not the question. But let us actually *build* another clock which is exactly the same as the first clock looks in the mirror—every time there is a screw with a right-hand thread in one, we use a screw with a left-hand thread in the corresponding place of the other; where one is marked “2” on the face, we mark a “5” on the face of the other; each coiled spring is twisted one way in one clock and the other way in the mirror-image clock; when we are all finished, we have two clocks, both physical, which bear to each other the relation of an object and its mirror image, although they are both actual, material objects, we emphasize. Now the question is: If the two clocks are started in the same condition, the springs wound to corresponding tightnesses, will the two clocks tick and go around, forever after, as exact mirror images? (This is a physical question, not a philosophical question.) Our intuition about the laws of physics would suggest that they *would*.

We would suspect that, at least in the case of these clocks, reflection in space is one of the symmetries of physical laws, that if we change everything from “right” to “left” and leave it otherwise the same, we cannot tell the difference. Let us, then, suppose for a moment that this is true. If it is true, then it would be impossible to distinguish “right” and “left” by any physical phenomenon, just as it is, for example, impossible to define a particular absolute velocity by a physical phenomenon. So it should be impossible, by any physical phenomenon, to define absolutely what we mean by “right” as opposed to “left,” because the physical laws should be symmetrical.

Of course, the world does not *have* to be symmetrical. For example, using what we may call “geography,” surely “right” can be defined. For instance, we stand

in New Orleans and look at Chicago, and Florida is to our right (when our feet are on the ground!). So we can define “right” and “left” by geography. Of course, the actual situation in any system does not have to have the symmetry that we are talking about; it is a question of whether the *laws* are symmetrical—in other words, whether it is *against the physical laws* to have a sphere like the earth with “left-handed dirt” on it and a person like ourselves standing looking at a city like Chicago from a place like New Orleans, but with everything the other way around, so Florida is on the other side. It clearly seems not impossible, not against the physical laws, to have everything changed left for right.

Another point is that our definition of “right” should not depend on history. An easy way to distinguish right from left is to go to a machine shop and pick up a screw at random. The odds are it has a right-hand thread—not necessarily, but it is much more likely to have a right-hand thread than a left-hand one. This is a question of history or convention, or the way things happen to be, and is again not a question of fundamental laws. As we can well appreciate, everyone could have started out making left-handed screws!

So we must try to find some phenomenon in which “right hand” is involved fundamentally. The next possibility we discuss is the fact that polarized light rotates its plane of polarization as it goes through, say, sugar water. As we saw in Chapter 33, it rotates, let us say, to the right in a certain sugar solution. That is a way of defining “right-hand,” because we may dissolve some sugar in the water and then the polarization goes to the right. But sugar has come from living things, and if we try to make the sugar artificially, then we discover that it *does not* rotate the plane of polarization! But if we then take that same sugar which is made artificially and which does not rotate the plane of polarization, and put bacteria in it (they eat some of the sugar) and then filter out the bacteria, we find that we still have sugar left (almost half as much as we had before), and this time it does rotate the plane of polarization, but *the other way*! It seems very confusing, but is easily explained.

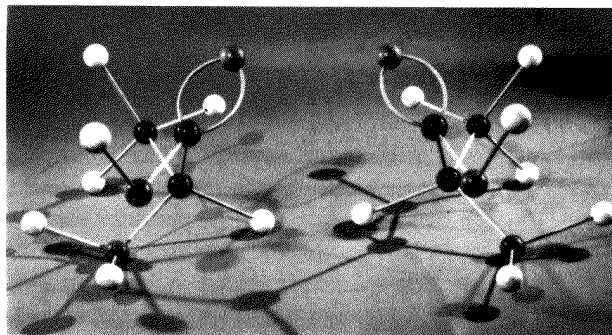


Fig. 52-1. (a) L-alanine (left), and (b) D-alanine (right).

Take another example: One of the substances which is common to all living creatures and that is fundamental to life is protein. Proteins consist of chains of amino acids. Figure 52-1 shows a model of an amino acid that comes out of a protein. This amino acid is called alanine, and the molecular arrangement would look like that in Fig. 52-1(a) if it came out of a protein of a real living thing. On the other hand, if we try to make alanine from carbon dioxide, ethane, and ammonia (and we *can* make it, it is not a complicated molecule), we discover that we are making equal amounts of this molecule and the one shown in Fig. 52-1(b)! The first molecule, the one that comes from the living thing, is called *L-alanine*. The other one, which is the same chemically, in that it has the same kinds of atoms and the same connections of the atoms, is a “right-hand” molecule, compared with the “left-hand” L-alanine, and it is called *D-alanine*. The interesting thing is that when we make alanine at home in a laboratory from simple gases, we get an equal mixture of both kinds. However, the only thing that life uses is L-alanine. (This is not exactly true. Here and there in living creatures there is a special use for D-alanine, but it is very rare. All proteins use L-alanine exclusively.) Now if we make both kinds, and we feed the mixture to some animal which likes to “eat,” or use up, alanine, it cannot use D-alanine, so it only uses the L-alanine; that is what happened to our sugar—after the bacteria eat the sugar that works

well for them, only the “wrong” kind is left! (Left-handed sugar tastes sweet, but not the same as right-handed sugar.)

So it looks as though the phenomena of life permit a distinction between “right” and “left,” or chemistry permits a distinction, because the two molecules are chemically different. But no, it does not! So far as physical measurements can be made, such as of energy, the rates of chemical reactions, and so on, the two kinds work exactly the same way if we make everything else in a mirror image too. One molecule will rotate light to the right, and the other will rotate it to the left in precisely the same amount, through the same amount of fluid. Thus, so far as physics is concerned, these two amino acids are equally satisfactory. So far as we understand things today, the fundamentals of the Schrödinger equation have it that the two molecules should behave in exactly corresponding ways, so that one is to the right as the other is to the left. Nevertheless, in life it is all one way!

It is presumed that the reason for this is the following. Let us suppose, for example, that life is somehow at one moment in a certain condition in which all the proteins in some creatures have left-handed amino acids, and all the enzymes are lopsided—every substance in the living creature is lopsided—it is not symmetrical. So when the digestive enzymes try to change the chemicals in the food from one kind to another, one kind of chemical “fits” into the enzyme, but the other kind does not (like Cinderella and the slipper, except that it is a “left foot” that we are testing). So far as we know, in principle, we could build a frog, for example, in which every molecule is reversed, everything is like the “left-hand” mirror image of a real frog; we have a left-hand frog. This left-hand frog would go on all right for a while, but he would find nothing to eat, because if he swallows a fly, his enzymes are not built to digest it. The fly has the wrong “kind” of amino acids (unless we give him a left-hand fly). So as far as we know, the chemical and life processes would continue in the same manner if everything were reversed.

If life is entirely a physical and chemical phenomenon, then we can understand that the proteins are all made in the same corkscrew only from the idea that at the very beginning some living molecules, by accident, got started and a few won. Somewhere, once, one organic molecule was lopsided in a certain way, and from this particular thing the “right” happened to evolve in our particular geography; a particular historical accident was one-sided, and ever since then the lopsidedness has propagated itself. Once having arrived at the state that it is in now, of course, it will always continue—all the enzymes digest the right things, manufacture the right things: when the carbon dioxide and the water vapor, and so on, go in the plant leaves, the enzymes that make the sugars make them lopsided because the enzymes are lopsided. If any new kind of virus or living thing were to originate at a later time, it would survive only if it could “eat” the kind of living matter already present. Thus it, too, must be of the same kind.

There is no conservation of the number of right-handed molecules. Once started, we could keep increasing the number of right-handed molecules. So the presumption is, then, that the phenomena in the case of life do not show a lack of symmetry in physical laws, but do show, on the contrary, the universal nature and the commonness of ultimate origin of all creatures on earth, in the sense described above.

52-5 Polar and axial vectors

Now we go further. We observe that in physics there are a lot of other places where we have “right” and “left” hand rules. As a matter of fact, when we learned about vector analysis we learned about the right-hand rules we have to use in order to get the angular momentum, torque, magnetic field, and so on, to come out right. The force on a charge moving in a magnetic field, for example, is $\mathbf{F} = q\mathbf{v} \times \mathbf{B}$. In a given situation, in which we know \mathbf{F} , \mathbf{v} , and \mathbf{B} , isn't that equation enough to define right-handedness? As a matter of fact, if we go back and look at where the vectors came from, we know that the “right-hand rule” was merely a convention; it was a trick. The original quantities, like the angular momenta and the angular velocities, and things of this kind, were not really vectors at all! They are all

somehow associated with a certain plane, and it is just because there are three dimensions in space that we can associate the quantity with a direction perpendicular to that plane. Of the two possible directions, we chose the “right-hand” direction.

So if the laws of physics are symmetrical, we should find that if some demon were to sneak into all the physics laboratories and replace the word “right” for “left” in every book in which “right-hand rules” are given, and instead we were to use all “left-hand rules,” uniformly, then it should make no difference whatever in the physical laws.

Let us give an illustration. There are two kinds of vectors. There are “honest” vectors, for example a step $\Delta \mathbf{r}$ in space. If in our apparatus there is a piece here and something else there, then in a mirror apparatus there will be the image piece and the image something else, and if we draw a vector from the “piece” to the “something else,” one vector is the mirror image of the other (Fig. 52-2). The vector arrow changes its head, just as the whole space turns inside out; such a vector we call a *polar vector*.

But the other kind of vector, which has to do with rotations, is of a different nature. For example, suppose that in three dimensions something is rotating as shown in Fig. 52-3. Then if we look at it in a mirror, it will be rotating as indicated, namely, as the mirror image of the original rotation. Now we have agreed to represent the mirror rotation by the same rule, it is a “vector” which, on reflection, does *not* change about as the polar vector does, but is reversed relative to the polar vectors and to the geometry of the space; such a vector is called an *axial vector*.

Now if the law of reflection symmetry is right in physics, then it must be true that the equations must be so designed that if we change the sign of each axial vector and each cross-product of vectors, which would be what corresponds to reflection, nothing will happen. For instance, when we write a formula which says that the angular momentum is $\mathbf{L} = \mathbf{r} \times \mathbf{p}$, that equation is all right, because if we change to a left-hand coordinate system, we change the sign of \mathbf{L} , but \mathbf{p} and \mathbf{r} do not change; the cross-product sign is changed, since we must change from a right-hand rule to a left-hand rule. As another example, we know that the force on a charge moving in a magnetic field is $\mathbf{F} = q\mathbf{v} \times \mathbf{B}$, but if we change from a right- to a left-handed system, since \mathbf{F} and \mathbf{v} are known to be polar vectors the sign change required by the cross-product must be cancelled by a sign change in \mathbf{B} , which means that \mathbf{B} must be an axial vector. In other words, if we make such a reflection, \mathbf{B} must go to $-\mathbf{B}$. So if we change our coordinates from right to left, we must also change the poles of magnets from north to south.

Let us see how that works in an example. Suppose that we have two magnets, as in Fig. 52-4. One is a magnet with the coils going around a certain way, and with current in a given direction. The other magnet looks like the reflection of the first magnet in a mirror—the coil will wind the other way, everything that happens inside the coil is exactly reversed, and the current goes as shown. Now, from the laws for the production of magnetic fields, which we do not know yet officially, but which we most likely learned in high school, it turns out that the magnetic field is as shown in the figure. In one case the pole is a south magnetic pole, while in the other magnet the current is going the other way and the magnetic field is reversed—it is a north magnetic pole. So we see that when we go from right to left we must indeed change from north to south!

Never mind changing north to south; these too are mere conventions. Let us talk about *phenomena*. Suppose, now, that we have an electron moving through one field, going into the page. Then, if we use the formula for the force, $\mathbf{v} \times \mathbf{B}$ (remember the charge is minus), we find that the electron will deviate in the indicated direction according to the physical law. So the phenomenon is that we have a coil with a current going in a specified sense and an electron curves in a certain way—that is the physics—never mind how we label everything.

Now let us do the same experiment with a mirror: we send an electron through in a corresponding direction and now the force is reversed, if we calculate it from the same rule, and that is very good because the corresponding *motions* are then mirror images!



Fig. 52-2. A step in space and its mirror image.

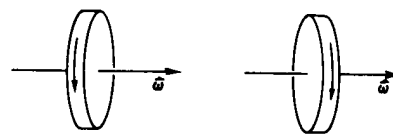


Fig. 52-3. A rotating wheel and its mirror image. Note that the angular velocity “vector” is not reversed in direction.

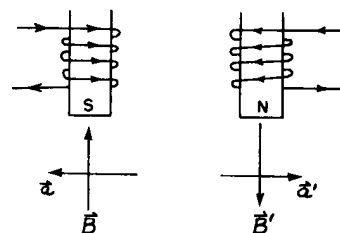


Fig. 52-4. A magnet and its mirror image.

52-6 Which hand is right?

So the fact of the matter is that in studying any phenomenon there are always two right-hand rules, or an even number of them, and the net result is that the phenomena always look symmetrical. In short, therefore, we cannot tell right from left if we also are not able to tell north from south. However, it may seem that we *can* tell the north pole of a magnet. The north pole of a compass needle, for example, is one that points to the north. But of course that is again a local property that has to do with geography of the earth; that is just like talking about in which direction is Chicago, so it does not count. If we have seen compass needles, we may have noticed that the north-seeking pole is a sort of bluish color. But that is just due to the man who painted the magnet. These are all local, conventional criteria.

However, if a magnet were to have the property that if we looked at it closely enough we would see small hairs growing on its north pole but not on its south pole, if that were the general rule, or if there were *any* unique way to distinguish the north from the south pole of a magnet, then we could tell which of the two cases we actually had, and *that would be the end of the law of reflection symmetry*.

To illustrate the whole problem still more clearly, imagine that we were talking to a Martian, or someone very far away, by telephone. We are not allowed to send him any actual samples to inspect; for instance, if we could send light, we could send him right-hand circularly polarized light and say, "That is right-hand light—just watch the way it is going." But we cannot *give* him anything, we can only talk to him. He is far away, or in some strange location, and he cannot see anything we can see. For instance, we cannot say, "Look at Ursa major; now see how those stars are arranged. What we mean by 'right' is . . ." We are only allowed to telephone him.

Now we want to tell him all about us. Of course, first we start defining numbers, and say, "Tick, tick, *two*, tick, tick, tick, *three* . . .," so that gradually he can understand a couple of words, and so on. After a while we may become very familiar with this fellow, and he says, "What do you guys look like?" We start to describe ourselves, and say, "Well, we are six feet tall." He says, "Wait a minute, what is six feet?" Is it possible to tell him what six feet is? Certainly! We say, "You know about the diameter of hydrogen atoms—we are 17,000,000,000 hydrogen atoms high!" That is possible because physical laws are not variant under change of scale, and therefore we *can* define an absolute length. And so we define the size of the body, and tell him what the general shape is—it has prongs with five bumps sticking out on the ends, and so on, and he follows us along, and we finish describing how we look on the outside, presumably without encountering any particular difficulties. He is even making a model of us as we go along. He says, "My, you are certainly very handsome fellows; now what is on the inside?" So we start to describe the various organs on the inside, and we come to the heart, and we carefully describe the shape of it, and say, "Now put the heart on the left side." He says, "Duhhh—the left side?" Now our problem is to describe to him which side the heart goes on without his ever seeing anything that we see, and without our ever sending any sample to him of what we mean by "right"—no standard right-handed object. Can we do it?

52-7 Parity is not conserved!

It turns out that the laws of gravitation, the laws of electricity and magnetism, nuclear forces, all satisfy the principle of reflection symmetry, so these laws, or anything derived from them, cannot be used. But associated with the many particles that are found in nature there is a phenomenon called *beta decay*, or *weak decay*. One of the examples of weak decay, in connection with a particle discovered in about 1954, posed a strange puzzle. There was a certain charged particle which disintegrated into three π -mesons, as shown schematically in Fig. 52-5. This particle was called, for a while, a τ -meson. Now in Fig. 52-5 we also see another particle which disintegrates into *two* mesons; one must be neutral, from the con-

servation of charge. This particle was called a θ -meson. So on the one hand we have a particle called a τ , which disintegrates into three π -mesons, and a θ , which disintegrates into two π -mesons. Now it was soon discovered that the τ and the θ are almost equal in mass; in fact, within the experimental error, they are equal. Next, the length of time it took for them to disintegrate into three π 's and two π 's was found to be almost exactly the same; they live the same length of time. Next, whenever they were made, they were made in the same proportions, say, 14 percent τ 's to 86 percent θ 's.

Anyone in his right mind realizes immediately that they must be the same particle, that we merely produce an object which has two different ways of disintegrating—not two different particles. This object that can disintegrate in two different ways has, therefore, the same lifetime and the same production ratio (because this is simply the ratio of the odds with which it disintegrates into these two kinds).

However, it was possible to prove (and we cannot here explain at all *how*), from the principle of reflection symmetry in quantum mechanics, that it was *impossible* to have these both come from the same particle—the same particle *could not* disintegrate in both of these ways. The conservation law corresponding to the principle of reflection symmetry is something which has no classical analog, and so this kind of quantum-mechanical conservation was called the *conservation of parity*. So, it was a result of the conservation of parity or, more precisely, from the symmetry of the quantum-mechanical equations of the weak decays under reflection, that the same particle could not go into both, so it must be some kind of coincidence of masses, lifetimes, and so on. But the more it was studied, the more remarkable the coincidence, and the suspicion gradually grew that possibly the deep law of the reflection symmetry of nature may be false.

As a result of this apparent failure, the physicists Lee and Yang suggested that other experiments be done in related decays to try to test whether the law was correct in other cases. The first such experiment was carried out by Miss Wu from Columbia, and was done as follows. Using a very strong magnet at a very low temperature, it turns out that a certain isotope of cobalt, which disintegrates by emitting an electron, is magnetic, and if the temperature is low enough that the thermal oscillations do not jiggle the atomic magnets about too much, they line up in the magnetic field. So the cobalt atoms will all line up in this strong field. They then disintegrate, emitting an electron, and it was discovered that when the atoms were lined up in a field whose \mathbf{B} vector points upward, most of the electrons were emitted in a downward direction.

If one is not really “hep” to the world, such a remark does not sound like anything of significance, but if one appreciates the problems and interesting things in the world, then he sees that it is a most dramatic discovery: When we put cobalt atoms in an extremely strong magnetic field, more disintegration electrons go down than up. Therefore if we were to put it in a corresponding experiment in a “mirror,” in which the cobalt atoms would be lined up in the opposite direction, they would spit their electrons *up*, not *down*; the action is *unsymmetrical*. *The magnet has grown hairs!* The south pole of a magnet is of such a kind that the electrons in a β -disintegration tend to go away from it; that distinguishes, in a physical way, the north pole from the south pole.

After this, a lot of other experiments were done: the disintegration of the π into μ and ν ; μ into an electron and two neutrinos; nowadays, the Λ into proton and π ; disintegration of Σ 's; and many other disintegrations. In fact, in almost all cases where it could be expected, all have been found *not* to obey reflection symmetry! Fundamentally, the law of reflection symmetry, at this level in physics, is incorrect.

In short, we can tell a Martian where to put the heart: we say, “Listen, build yourself a magnet, and put the coils in, and put the current on, and then take some cobalt and lower the temperature. Arrange the experiment so the electrons go from the foot to the head, then the direction in which the current goes through the coils is the direction that goes in on what we call the right and comes out on the left.” So it is possible to define right and left, now, by doing an experiment of this kind.

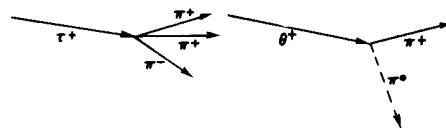


Fig. 52-5. A schematic diagram of the disintegration of a τ^+ and a θ^+ particle.

There are a lot of other features that were predicted. For example, it turns out that the spin, the angular momentum, of the cobalt nucleus before disintegration is 5 units of \hbar , and after disintegration it is 4 units. The electron carries spin angular momentum, and there is also a neutrino involved. It is easy to see from this that the electron must carry its spin angular momentum aligned along its direction of motion, the neutrino likewise. So it looks as though the electron is spinning to the left, and that was also checked. In fact, it was checked right here at Caltech by Boehm and Wapstra, that the electrons spin mostly to the left. (There were some other experiments that gave the opposite answer, but they were wrong!)

The next problem, of course, was to find the law of the failure of parity conservation. What is the rule that tells us how strong the failure is going to be? The rule is this: it occurs only in these very slow reactions, called weak decays, and when it occurs, the rule is that the particles which carry spin, like the electron, neutrino, and so on, come out with a spin tending to the left. That is a lopsided rule; it connects a polar vector velocity and an axial vector angular momentum, and says that the angular momentum is more likely to be opposite to the velocity than along it.

Now that is the rule, but today we do not really understand the whys and wherefores of it. *Why* is this the right rule, what is the fundamental reason for it, and how is it connected to anything else? At the moment we have been so shocked by the fact that this thing is unsymmetrical that we have not been able to recover enough to understand what it means with regard to all the other rules. However, the subject is interesting, modern, and still unsolved, so it seems appropriate that we discuss some of the questions associated with it.

52-8 Antimatter

The first thing to do when one of the symmetries is lost is to immediately go back over the list of known or assumed symmetries and ask whether any of the others are lost. Now we did not mention one operation on our list, which must necessarily be questioned, and that is the relation between matter and antimatter. Dirac predicted that in addition to electrons there must be another particle, called the positron (discovered at Caltech by Anderson), that is necessarily related to the electron. All the properties of these two particles obey certain rules of correspondence: the energies are equal; the masses are equal; the charges are reversed; but, more important than anything, the two of them, when they come together, can annihilate each other and liberate their entire mass in the form of energy, say γ -rays. The positron is called an *antiparticle* to the electron, and these are the characteristics of a particle and its antiparticle. It was clear from Dirac's argument that all the rest of the particles in the world should also have corresponding antiparticles. For instance, for the proton there should be an antiproton, which is now symbolized by a \bar{p} . The \bar{p} would have a negative electrical charge and the same mass as a proton, and so on. The most important feature, however, is that a proton and an antiproton coming together can annihilate each other. The reason we emphasize this is that people do not understand it when we say there is a neutron and also an antineutron, because they say, "A neutron is neutral, so how *can* it have the opposite charge?" The rule of the "anti" is not just that it has the opposite charge, it has a certain set of properties, the whole lot of which are opposite. The antineutron is distinguished from the neutron in this way: if we bring two neutrons together, they just stay as two neutrons, but if we bring a neutron and an antineutron together, they annihilate each other with a great explosion of energy being liberated, with various π -mesons, γ -rays, and whatnot.

Now if we have antineutrons, antiprotons, and antielectrons, we can make antiatoms, in principle. They have not been made yet, but it is possible in principle. For instance, a hydrogen atom has a proton in the center with an electron going around outside. Now imagine that somewhere we can make an antiproton with a positron going around, would it go around? Well, first of all, the antiproton is electrically negative and the antielectron is electrically positive, so they attract each other in a corresponding manner—the masses are all the same; everything is the

same. It is one of the principles of the symmetry of physics, the equations seem to show, that if a clock, say, were made of matter on one hand, and then we made the same clock of antimatter, it would run in this way. (Of course, if we put the clocks together, they would annihilate each other, but that is different.)

An immediate question then arises. We can build, out of matter, two clocks, one which is "left-hand" and one which is "right-hand." For example, we could build a clock which is not built in a simple way, but has cobalt and magnets and electron detectors which detect the presence of β -decay electrons and count them. Each time one is counted, the second hand moves over. Then the mirror clock, receiving fewer electrons, will not run at the same rate. So evidently we can make two clocks such that the left-hand clock does not agree with the right-hand one. Let us make, out of matter, a clock which we call the standard or right-hand clock. Now let us make, also out of matter, a clock which we call the left-hand clock. We have just discovered that, in general, these two will *not* run the same way; prior to that famous physical discovery, it was thought that they would. Now it was also supposed that matter and antimatter were equivalent. That is, if we made an antimatter clock, right-hand, the same shape, then it would run the same as the right-hand matter clock, and if we made the same clock to the left it would run the same. In other words, in the beginning it was believed that *all four* of these clocks were the same; now of course we know that the right-hand and left-hand matter are not the same. Presumably, therefore, the right-handed antimatter and the left-handed antimatter are not the same.

So the obvious question is, which goes with which, if either? In other words, does the right-handed matter behave the same way as the right-handed antimatter? Or does the right-handed matter behave the same as the left-handed antimatter? β -decay experiments, using positron decay instead of electron decay, indicate that this is the interconnection: matter to the "right" works the same way as antimatter to the "left."

Therefore, at long last, it is really true that right and left symmetry is still maintained! If we made a left-hand clock, but made it out of the other kind of matter, antimatter instead of matter, it would run in the same way. So what has happened is that instead of having two independent rules in our list of symmetries, two of these rules go together to make a new rule, which says that matter to the right is symmetrical with antimatter to the left.

So if our Martian is made of antimatter and we give him instructions to make this "right" handed model like us, it will, of course, come out the other way around. What would happen when, after much conversation back and forth, we each have taught the other to make space ships and we meet halfway in empty space? We have instructed each other on our traditions, and so forth, and the two of us come rushing out to shake hands. Well, if he puts out his left hand, watch out!

52-9 Broken symmetries

The next question is, what can we make out of laws which are *nearly* symmetrical? The marvelous thing about it all is that for such a wide range of important, strong phenomena—nuclear forces, electrical phenomena, and even weak ones like gravitation—over a tremendous range of physics, all the laws for these seem to be symmetrical. On the other hand, this little extra piece says, "No, the laws are not symmetrical!" How is it that nature can be almost symmetrical, but not perfectly symmetrical? What shall we make of this? First, do we have any other examples? The answer is, we do, in fact, have a few other examples. For instance, the nuclear part of the force between proton and proton, between neutron and neutron, and between neutron and proton, is all exactly the same—there is a symmetry for nuclear forces, a new one, that we can interchange neutron and proton—but it evidently is not a general symmetry, for the electrical repulsion between two protons at a distance does not exist for neutrons. So it is not generally true that we can *always* replace a proton with a neutron, but only to a good approximation. Why *good*? Because the nuclear forces are much stronger than the

electrical forces. So this is an “almost” symmetry also. So we do have examples in other things.

We have, in our minds, a tendency to accept symmetry as some kind of perfection. In fact it is like the old idea of the Greeks that circles were perfect, and it was rather horrible to believe that the planetary orbits were not circles, but only nearly circles. The difference between being a circle and being nearly a circle is not a small difference, it is a fundamental change so far as the mind is concerned. There is a sign of perfection and symmetry in a circle that is not there the moment the circle is slightly off—that is the end of it—it is no longer symmetrical. Then the question is why it is only *nearly* a circle—that is a much more difficult question. The actual motion of the planets, in general, should be ellipses, but during the ages, because of tidal forces, and so on, they have been made almost symmetrical. Now the question is whether we have a similar problem here. The problem from the point of view of the circles is if they were perfect circles there would be nothing to explain, that is clearly simple. But since they are only nearly circles, there is a lot to explain, and the result turned out to be a big dynamical problem, and now our problem is to explain why they are nearly symmetrical by looking at tidal forces and so on.

So our problem is to explain where symmetry comes from. Why is nature so nearly symmetrical? No one has any idea why. The only thing we might suggest is something like this: There is a gate in Japan, a gate in Neiko, which is sometimes called by the Japanese the most beautiful gate in all Japan; it was built in a time when there was great influence from Chinese art. This gate is very elaborate, with lots of gables and beautiful carving and lots of columns and dragon heads and princes carved into the pillars, and so on. But when one looks closely he sees that in the elaborate and complex design along one of the pillars, one of the small design elements is carved upside down; otherwise the thing is completely symmetrical. If one asks why this is, the story is that it was carved upside down so that the gods will not be jealous of the perfection of man. So they purposely put an error in there, so that the gods would not be jealous and get angry with human beings.

We might like to turn the idea around and think that the true explanation of the near symmetry of nature is this: that God made the laws only nearly symmetrical so that we should not be jealous of His perfection!

Index

- Aberration, 27–7, 34–10
- Absolute zero, 1–5
- Absorption, 31–8 ff
- Acceleration, 8–8 ff
 - components of, 9–3
 - of gravity, 9–4
- Activation energy, 42–7
- Adams, J. C., 7–5
- Adiabatic compression, 39–5
- Adiabatic expansion, 44–5
- Affective future, 17–4
- Air trough, 10–5
- Algebra, 22–1 ff
- Amplitude modulation, 48–3
- Amplitude, of oscillation, 21–3
- Analog computer, 25–8
- Anderson, C. D., 52–10
- Angle, of incidence, 26–3
 - of reflection, 26–3
- Angstrom (unit), 1–3
- Angular frequency, 21–3, 29–2
- Angular momentum, 7–7, 18–5 f, 20–1
 - conservation of, 4–7, 18–6 ff, 20–5
 - of rigid body, 20–8
- Anomalous refraction, 33–9 f
- Antimatter 52–10 f
- Antiparticle, 2–8
- Aristotle, 5–1
- Atom, 1–2
 - metastable, 42–10
- Atomic clock, 5–5
- Atomic hypothesis, 1–2
- Atomic particles, 2–9 f
- Atomic processes, 1–5 f
- Attenuation, 31–8
- Avogadro, A., 39–2
- Avogadro's number, 41–10
- Axial vector, 52–6 f

- Becquerel, A. H., 28–3
- Birefringence, 33–3 ff
- Blackbody radiation, 41–5 f
- Bohm, 52–10
- Bohr, N., 42–9
- Bohr radius, 38–6
- Boltzmann, 41–2
- Boltzmann's law, 40–2 f
- Born, M., 37–1, 38–9
- Boyle's law, 40–8
- Bremsstrahlung, 34–6 f
- Brewster's angle, 33–6
- Briggs, H., 22–6
- Brown, R., 41–1
- Brownian motion, 1–8, 6–5, 41–1 ff

- Capacitance, 23–5
- Capacitor, 14–9, 23–5
- Capillary action, 51–8
- Carnot, S., 4–2, 44–2 ff
- Carnot cycle, 44–5 f, 45–2
- Carrier signal, 48–3
- Catalyst, 42–8

- Cavendish, H., 7–9
- Cavendish's experiment, 7–9
- Center of mass, 18–1 f, 19–1 ff
- Centrifugal force, 7–5, 12–11
- Čerenkov, P. A., 51–2
- Čerenkov radiation, 51–2
- Charge, conservation of, 4–7
 - on electron, 12–7
- Chemical energy, 4–2
- Chemical kinetics, 42–7 f
- Chemical reaction, 1–6 ff
- Chromaticity, 35–6 f
- Circular motion, 21–4
- Clausius, R., 44–2, 44–3
- Clausius-Clapeyron equation, 45–6 ff
- Coefficient, of friction, 12–4
 - gravitational, 7–9
- Collision, 16–6
 - elastic, 10–7
- Color vision, 35–1 ff
 - physiochemistry of, 35–9 f
- Complex impedance, 23–7
- Complex numbers, 22–7 ff, 23–1 ff
- Compound eye, 36–6 ff
- Compression, adiabatic, 39–5
 - isothermal, 44–5
- Cones, 35–1
- Conservation, of angular momentum, 4–7, 18–6 ff, 20–5
 - of charge, 4–7
 - of energy, 3–2, 4–1 ff
 - of linear momentum, 4–7, 10–1 ff
- Contraction hypothesis, 15–3
- Copernicus, 7–1
- Coriolis force, 19–8 f
- Cornea, 35–1
- Coulomb's law, 28–2
- Cross section, for scattering, 32–7
- Crystal diffraction, 38–4 f

- Dedekind, R., 22–4
- Degrees of freedom, 25–2, 39–12
- Density, 1–4
- Derivative, 8–5 ff
 - partial, 14–9
- Dicke, R. H., 7–11
- Differential calculus, 8–4
- Diffraction, 30–1 ff
 - by screen, 31–10 f
- Diffraction grating, 29–5, 30–3 ff
- Diffusion, 43–1 ff
- Dipole moment, 12–6
- Dipole radiator, 28–5 f, 29–3 ff
- Dirac, P., 52–10
- Dirac equation, 20–6
- Dispersion, 31–6 ff
- Distance, 5–5 ff
- Distance measurement, color brightness, 5–6
 - triangulation, 5–6
- Doppler effect, 17–8, 23–9, 34–7 f, 38–6
- Double stars, 7–6

- Dynamics, 7–2 f, 9–1 ff
 - relativistic, 15–9 f

- Efficiency, of ideal engine, 44–7 f
- Einstein, A., 2–6, 7–11, 12–12, 15–1, 16–1, 41–8, 42–8, 42–9
- Elastic collision, 10–7
- Elastic energy, 4–2, 4–6
- Electrical energy, 4–2
- Electric field, 2–4, 12–7 f
- Electromagnetic energy, 29–2
- Electromagnetic field, 2–2, 2–5, 10–9
- Electromagnetic radiation, 26–1, 28–1 ff
- Electromagnetic waves, cosmic rays, 2–5
 - gamma rays, 2–5
 - infrared, 2–5, 23–8, 26–1
 - light, 2–5
 - ultraviolet, 2–5, 26–1
 - x-rays, 2–5, 26–1
- Electron, 2–4, 37–1, 37–4 ff
 - charge on, 12–7
 - radius of, classical, 32–4
- Electron cloud, 6–11
- Electron-ray tube, 12–9
- Electron volt (unit), 34–4
- Ellipse, 7–1
- Energy, chemical, 4–2
 - conservation of, 3–2, 4–1 ff
 - elastic, 4–2, 4–6
 - electrical, 4–2
 - electromagnetic, 29–2
 - gravitational, 4–2 ff
 - heat, 4–2, 4–6, 10–7, 10–8
 - kinetic, 1–7, 4–2, 4–5 f, 39–4
 - mass, 4–2, 4–7
 - nuclear, 4–2
 - potential, 4–4, 13–1 ff, 14–1 ff
 - radiant, 4–2
 - relativistic, 16–1 ff
- Energy levels, 38–7 f
- Energy theorem, 50–7 f
- Enthalpy, 45–5
- Entropy, 44–10 ff, 46–7 ff
- Eötvös, 7–11
- Equilibrium, 1–6
- Euclid, 5–6
- Euclidean geometry, 12–3
- Evaporation, 1–5 f
 - of a liquid, 40–3 f, 42–1 ff
- Expansion, adiabatic, 44–5
 - isothermal, 44–5
- Exponential atmosphere, 40–1 f
- Eye, compound, 36–6 ff
 - human, 35–1 f, 36–3 ff

- Farad (unit), 25–7
- Fermat, P., 26–3
- Fermi (unit), 5–10
- Fermi, E., 5–10
- Fields, 2–2, 2–4, 2–5, 10–9, 12–7 ff, 13–8 f, 14–7 ff
 - superposition of, 12–9

- Focal length, 27-1 ff
Focus, 26-5
Force, centrifugal, 7-5, 12-11
 components of, 9-3
 conservative, 14-3 ff
 Coriolis, 19-8 f
 electrical, 2-3 ff
 gravitational, 2-3
 molecular, 1-3, 12-6 f
 moment of, 18-5
 nonconservative, 14-6 f
 nuclear, 12-12
 pseudo, 12-10 ff
Fourier, J., 50-2 f
Fourier analysis, 50-2 ff
Fourier transform, 25-4
Four-vectors, 15-8 f, 17-5 ff
Fovea, 35-1
Frank, I., 51-2
Frequency, angular, 21-3, 29-2
 of oscillation, 2-5
Fresnel's reflection formulas, 33-8
Friction, 10-5, 12-3 ff
 coefficient of, 12-4
- Galileo, 5-1, 7-2, 9-1, 52-3
Galilean relativity, 10-3
Galilean transformation, 12-11
Gauss (unit), 34-4
Gell-Mann, M., 2-9
Geometrical optics, 26-1, 27-1 f
Gravitation, 2-3, 7-1 ff, 12-2
Gravitational acceleration, 9-4
Gravitational coefficient, 7-9
Gravitational energy, 4-2 ff
Gravitational field, 12-8 ff, 13-8 f
Gravity, 13-3 ff
 acceleration of, 9-4
Green's function, 25-4
Gyroscope, 20-5 ff
- Harmonic motion, 21-4, 23-1 ff
Harmonic oscillator, 10-1, 21-1 ff
 forced, 21-5 f, 23-3 ff
Harmonics, 50-1 ff
Heat, 1-3, 13-3
Heat energy, 4-2, 4-6, 10-7, 10-8
Heat engines, 44-1 ff
Heisenberg, W., 6-10, 37-1, 37-9, 37-11, 37-12, 38-9
Helmholtz, H., 35-7
Henry (unit), 25-7
Hooke's law, 12-6
Huygens, C., 15-2, 26-2
Hypocycloid, 34-3
- Ideal gas law, 39-10 ff
Impedance, 25-8 f
 complex, 23-7
Incidence, angle of, 26-3
Inclined plane, 4-4
Index, of refraction, 31-1 ff
Inductance, 23-6
Inductor, 23-6
Inertia, 2-3, 7-11
 moment of, 18-7, 19-5 ff
 principle of, 9-1
Infrared radiation, 23-8, 26-1
Integral, 8-7 f
Interference, 28-6, 29-1 ff
Interfering waves, 37-4
- Interferometer, 15-5
Ion, 1-6
Ionic conductivity, 43-6 f
Ionization energy, 42-5
Isothermal atmosphere, 40-2
Isothermal compression, 44-5
Isothermal expansion, 44-5
Isotopes, 3-4 ff
- Jeans, J., 40-9, 41-6 f
Johnson noise, 41-2, 41-8
Joule (unit), 13-3
Joule heating, 24-2
- Kepler, J., 7-1
Kepler's laws, 7-1 f, 9-1, 18-6
Kerr cell, 33-5
Kinetic energy, 1-7, 4-2, 4-5 f, 39-4
 rotational, 19-7 ff
Kinetic theory, 42-1 ff
 of gases, 39-1 ff
Kirchhoff's laws, 25-9
- Laplace, P., 47-7
Laser, 32-6, 42-10
Least time, principle of, 26-3 ff, 26-8
Leibnitz, G. W., 8-4
Lens formula, 27-6
Leverrier, U., 7-5
Light, momentum of, 34-10 f
 polarized, 32-9
 scattering of, 32-5 ff
 speed of, 15-1
Light waves, 48-1
Linear momentum, conservation of, 4-7, 10-1 ff
Linear systems, 25-1 ff
Logarithms, 22-4
Lorentz, H. A., 15-3
Lorentz contraction, 15-7
Lorentz transformation, 15-3, 17-1, 34-8, 52-2
- Magnetic field, 12-9 f
Magnetic induction, 12-10
Magnetism, 2-4
Magnification, 27-5
Maser, 42-10
Mass, 9-1, 15-1
 center of, 18-1 f, 19-1 ff
 relativistic, 16-6 ff
Mass energy, 4-2, 4-7
Mass-energy equivalence, 15-10 f
Maxwell, J. C., 6-1, 6-9, 28-1, 40-8, 41-7, 46-5
Maxwell's equations, 15-2, 25-3, 47-7
Mayer, J. R., 3-2
Mean free path, 43-3 f
Mean square distance, 6-5, 41-9
Mendeléev, 2-9
Metastable atom, 42-10
Meter (unit), 5-10
Mev (unit), 2-9
Michelson-Morley experiment, 15-3 ff
Miller, W. C., 35-2
Minkowski, 17-8
Modes, 49-1 ff
Mössbauer, R., 23-9
Mole (unit), 39-10
Molecular attraction, 1-3, 12-6 f
Molecular diffusion, 43-7 ff
- Molecular motion, 41-1
Molecule, 1-3
Moment, dipole, 12-6
 of force, 18-5
 of inertia, 18-7, 19-5 ff
Momentum, 9-1 f, 38-2 ff
 angular, 7-7, 18-5 ff, 20-1, 20-5
 of light, 34-10 f
 linear, 4-7, 10-1 ff
 relativistic, 10-8 f, 16-1 ff
Monatomic gas, 39-5
Motion, 5-1, 8-1 ff
 circular, 21-4
 constrained, 14-3
 harmonic, 21-4, 23-1 ff
 parabolic, 8-10
 planetary, 7-1 ff, 9-6 f, 13-5
Music, 50-1
- Nernst heat theorem, 44-11
Neutrons, 2-4
Newton, I., 8-4, 15-1, 37-1
Newton-meter (unit), 13-3
Newton's laws, 2-6, 7-3 ff, 7-11, 9-1 ff, 10-1 ff, 11-7 f, 12-1, 39-2, 41-1, 46-1
Nishijima, 2-9
Nodes, 49-2
Noise, 50-1
Nuclear cross section, 5-9
Nuclear energy, 4-2
Nuclear forces, 12-12
Nucleus, 2-4, 2-8 ff
Numerical analysis, 9-6
Nutation, 20-7
- Ohm (unit), 25-7
Ohm's law, 25-7, 43-7
Optic axis, 33-3
Optic nerve, 35-2
Optics, 26-1 ff
 geometrical, 26-1, 27-1 ff
Oscillation, amplitude of, 21-3
 damped, 24-3 f
 frequency of, 2-5
 period of, 21-3
 periodic, 9-4
 phase of, 21-3
Oscillator, 5-2
 harmonic, 10-1, 21-1, 21-5 f, 23-3 ff
- Pappus, theorem of, 19-4
Parabolic antenna, 30-6 f
Parabolic motion, 8-10
Parallel-axis theorem, 19-6
Parallel-plate capacitor, 14-9
Paraxial rays, 27-2
Partial derivative, 14-9
Pascal's triangle, 6-4
Pendulum, 49-6 f
Pendulum clock, 5-2
Period, of oscillation, 21-3
Periodic time, 5-1 f
Perpetual motion, 46-2
Phase, of oscillation, 21-3
Phase shift, 21-3
Phase velocity, 48-6
Photon, 2-7, 26-1, 37-8
Physiochemistry, of color vision, 35-9 f
Planck, M., 41-6, 42-8, 42-9
Planck's constant, 5-10, 6-10, 17-8, 37-11
Planetary motion, 7-1 ff, 9-6 f, 13-5

- Poincaré, H., 15-3, 15-5, 16-1
 Polarization, 33-1 ff
 Polarized light, 32-9
 Potential energy, 4-4, 13-1 ff, 14-1 ff
 Power, 13-2
 Pressure, 1-3
 Probability, 6-1 ff
 Probability density, 6-8 f
 Probability distribution, 6-7 ff
 Proton, 2-4
 Pseudo force, 12-10 ff
 Ptolemy, 26-2
 Purkinje effect, 35-2
 Pythagoras, 50-1
- Quantum electrodynamics, 2-7, 28-3
 Quantum mechanics, 2-2, 2-6 ff, 6-10, 10-9, 37-1 ff, 38-1 ff
- Radiant energy, 4-2
 Radiation, infrared, 23-8, 26-1
 relativistic effects, 35-1 ff
 synchrotron, 34-3 ff, 34-6
 ultraviolet, 26-1
 Radiation damping, 32-3 f
 Radiation resistance, 32-1 ff
 Radioactive clock, 5-3 ff
 Radius, of electron, 32-4
 Ramsey, N., 5-5
 Random walk, 6-5 ff, 41-8 ff
 Ratchet and pawl machine, 46-1 ff
 Rayleigh's criterion, 30-6
 Rayleigh's law, 41-6
 Reciprocity principle, 30-7
 Rectification, 50-9
 Reflection, 26-2 f
 angle of, 26-3
 Refraction, 26-2 f
 anomalous, 33-9 f
 index of, 31-1 ff
 Relativistic dynamics, 15-9 f
 Relativistic energy, 16-1 ff
 Relativistic mass, 16-1 ff
 Relativistic momentum, 10-8 f, 16-1 ff
 Relativity, special theory of, 15-1 ff
 Galilean, 10-3
 theory of, 7-11, 17-1
 Resistance, 23-5
 Resistor, 23-5
 Resolving power, 27-7 f, 30-5 f
 Resonance, 23-1 ff
 electrical, 23-5 ff
 in nature, 23-7 ff
 Resonance interaction, 2-9
 Retarded time, 28-2
 Retina, 35-1
 Rigid body, 18-1
 angular momentum of, 20-8
 rotation of, 18-2 ff
 Ritz combination principle, 38-8
- Rods, 35-1, 36-6
 Roemer, O., 7-5
 Root-mean-square distance, 6-6
 Rotation, of axes, 11-3 f
 plane, 18-1
 of a rigid body, 18-2 ff
 in space, 20-1 ff
 in two dimensions, 18-1 ff
 Rushton, 35-9
 Rydberg (unit), 38-6
- Scalar, 11-5
 Scattering, of light, 32-5 ff
 Schrödinger, E., 35-6, 37-1, 38-9
 Scientific method, 2-1 f
 Screw jack, 4-5
 Second (unit), 5-5
 Seismograph, 51-5
 Shannon, C., 44-2
 Shear wave, 51-4
 Side bands, 48-4 f
 Simultaneity, 15-7 f
 Sinusoidal waves, 29-2 f
 Smoluchowski, 41-8
 Smooth muscle, 14-2
 Snell, W., 26-3
 Snell's law, 26-3, 31-2
 Sound, 2-3, 47-1 ff, 50-1
 speed of, 47-7 f
 Space, 8-2
 Space-time, 2-6, 17-1 ff
 Special theory of relativity, 15-1 ff
 Specific heat, 40-7 f, 45-2
 Speed, 8-2 ff, 9-2
 of light, 15-1
 of sound, 47-7 f
 Spontaneous emission, 42-9
 Standard deviation, 6-9
 Statistical fluctuations, 6-3 ff
 Statistical mechanics, 3-1, 40-1 ff
 Stevinus, S., 4-5
 "Strangeness" number, 2-9
 Striated muscle, 14-2
 Superposition, of fields, 12-9
 principle of, 25-2 ff
 Symmetry, 1-4, 11-1 ff
 of physical laws, 16-3, 52-1 ff
 Synchrotron, 2-5, 15-9, 34-3 ff, 34-6
- Tamm, I., 51-2
 Temperature, 39-6 ff
 Thermal conductivity, of a gas, 43-9 f
 Thermal equilibrium, 41-3 ff
 Thermal ionization, 42-5 ff
 Thermodynamics, 39-2, 45-1 ff
 laws of, 44-1 ff
 Thompson scattering cross section, 32-8
 Three-body problem, 10-1
 Tides, 7-4 f
 Time, 2-3, 5-1 ff, 8-1, 8-2
 retarded, 28-2
- standard of, 5-5
 transformation of, 15-5 ff
 Torque, 18-4, 20-1 ff
 Transformation, Fourier, 25-4
 Galilean, 12-11
 linear, 11-6
 Lorentz, 15-3, 17-1, 34-8, 52-2
 of time, 15-5 ff
 of velocity, 16-4 ff
 Transient, 24-1 ff
 electrical, 24-5 f
 Transient response, 21-6
 Translation, of axes, 11-1 ff
 Twin paradox, 16-3 f
 Tycho Brahe, 7-1
- Ultraviolet radiation, 26-1
 Uncertainty principle, 2-6, 6-10 f, 37-9, 37-11, 38-8 f
 Unit cell, 38-5
 Unit vector, 11-10
- Vector, 11-5 ff
 Vector algebra, 11-6 f
 Vector analysis, 11-5, 52-2
 Vector product, 20-4
 Velocity, 8-3, 9-2 f
 components of, 9-3
 transformation of 16-4 ff
 Vinci, Leonardo da, 36-2
 Virtual work, principle of, 4-5
 Vision, 36-1 ff
 binocular, 36-4
 color, 35-1 ff
 Visual cortex, 36-4
 Visual purple, 35-9
- Wapstra, 52-10
 Watt (unit), 13-3
 Wave, shear, 51-4
 sinusoidal, 29-2 f
 Wave equation, 47-1 ff
 Wavefront, 47-3
 Wavelength, 19-3, 26-1
 Wave number, 29-2
 Waves, 51-1 ff
 light, 48-1
 Weyl, H., 11-1
 Work, 13-1 ff, 14-1 ff
- X-rays, 2-5, 26-1
- Young, 35-7
 Yukawa, H., 2-8
 Yustova, 35-8
- Zeno, 8-3
 Zero, absolute, 1-5
 Zero mass, 2-10

The Feynman

LECTURES ON PHYSICS

MAINLY ELECTROMAGNETISM AND MATTER

RICHARD P. FEYNMAN

*Richard Chace Tolman Professor of Theoretical Physics
California Institute of Technology*

ROBERT B. LEIGHTON

*Professor of Physics
California Institute of Technology*

MATTHEW SANDS

*Professor
Stanford University*

OXNARD PUBLIC LIBRARY
251 SOUTH A STREET
OXNARD, CALIFORNIA 93030



ADDISON-WESLEY PUBLISHING COMPANY, INC.

READING, MASSACHUSETTS • PALO ALTO • LONDON

Copyright © 1964

CALIFORNIA INSTITUTE OF TECHNOLOGY

Printed in the United States of America

ALL RIGHTS RESERVED. THIS BOOK, OR PARTS THEREOF
MAY NOT BE REPRODUCED IN ANY FORM WITHOUT WRITTEN
PERMISSION OF THE PUBLISHER

Library of Congress Catalog Card No. 63-20717

Second printing—November, 1964



Feynman's Preface

These are the lectures in physics that I gave last year and the year before to the freshman and sophomore classes at Caltech. The lectures are, of course, not verbatim—they have been edited, sometimes extensively and sometimes less so. The lectures form only part of the complete course. The whole group of 180 students gathered in a big lecture room twice a week to hear these lectures and then they broke up into small groups of 15 to 20 students in recitation sections under the guidance of a teaching assistant. In addition, there was a laboratory session once a week.

The special problem we tried to get at with these lectures was to maintain the interest of the very enthusiastic and rather smart students coming out of the high schools and into Caltech. They have heard a lot about how interesting and exciting physics is—the theory of relativity, quantum mechanics, and other modern ideas. By the end of two years of our previous course, many would be very discouraged because there were really very few grand, new, modern ideas presented to them. They were made to study inclined planes, electrostatics, and so forth, and after two years it was quite stultifying. The problem was whether or not we could make a course which would save the more advanced and excited student by maintaining his enthusiasm.

The lectures here are not in any way meant to be a survey course, but are very serious. I thought to address them to the most intelligent in the class and to make sure, if possible, that even the most intelligent student was unable to completely encompass everything that was in the lectures—by putting in suggestions of applications of the ideas and concepts in various directions outside the main line of attack. For this reason, though, I tried very hard to make all the statements as accurate as possible, to point out in every case where the equations and ideas fitted into the body of physics, and how—when they learned more—things would be modified. I also felt that for such students it is important to indicate what it is that they should—if they are sufficiently clever—be able to understand by deduction from what has been said before, and what is being put in as something new. When new ideas came in, I would try either to deduce them if they were deducible, or to explain that it *was* a new idea which hadn't any basis in terms of things they had already learned and which was not supposed to be provable—but was just added in.

At the start of these lectures, I assumed that the students knew something when they came out of high school—such things as geometrical optics, simple chemistry ideas, and so on. I also didn't see that there was any reason to make the lectures

in a definite order, in the sense that I would not be allowed to mention something until I was ready to discuss it in detail. There was a great deal of mention of things to come, without complete discussions. These more complete discussions would come later when the preparation became more advanced. Examples are the discussions of inductance, and of energy levels, which are at first brought in in a very qualitative way and are later developed more completely.

At the same time that I was aiming at the more active student, I also wanted to take care of the fellow for whom the extra fireworks and side applications are merely disquieting and who cannot be expected to learn most of the material in the lecture at all. For such students I wanted there to be at least a central core or backbone of material which he *could* get. Even if he didn't understand everything in a lecture, I hoped he wouldn't get nervous. I didn't expect him to understand everything, but only the central and most direct features. It takes, of course, a certain intelligence on his part to see which are the central theorems and central ideas, and which are the more advanced side issues and applications which he may understand only in later years.

In giving these lectures there was one serious difficulty: in the way the course was given, there wasn't any feedback from the students to the lecturer to indicate how well the lectures were going over. This is indeed a very serious difficulty, and I don't know how good the lectures really are. The whole thing was essentially an experiment. And if I did it again I wouldn't do it the same way—I hope I *don't* have to do it again! I think, though, that things worked out—so far as the physics is concerned—quite satisfactorily in the first year.

In the second year I was not so satisfied. In the first part of the course, dealing with electricity and magnetism, I couldn't think of any really unique or different way of doing it—of any way that would be particularly more exciting than the usual way of presenting it. So I don't think I did very much in the lectures on electricity and magnetism. At the end of the second year I had originally intended to go on, after the electricity and magnetism, by giving some more lectures on the properties of materials, but mainly to take up things like fundamental modes, solutions of the diffusion equation, vibrating systems, orthogonal functions, . . . developing the first stages of what are usually called “the mathematical methods of physics.” In retrospect, I think that if I were doing it again I would go back to that original idea. But since it was not planned that I would be giving these lectures again, it was suggested that it might be a good idea to try to give an introduction to the quantum mechanics—what you will find in Volume III.

It is perfectly clear that students who will major in physics can wait until their third year for quantum mechanics. On the other hand, the argument was made that many of the students in our course study physics as a background for their primary interest in other fields. And the usual way of dealing with quantum mechanics makes that subject almost unavailable for the great majority of students because they have to take so long to learn it. Yet, in its real applications—especially in its more complex applications, such as in electrical engineering and chemistry—the full machinery of the differential equation approach is not actually used. So I tried to describe the principles of quantum mechanics in a way which wouldn't require that one first know the mathematics of partial differential equations. Even for a physicist I think that is an interesting thing to try to do—to present quantum mechanics in this reverse fashion—for several reasons which may be apparent in the lectures themselves. However, I think that the experiment in the quantum mechanics part was not completely successful—in large part because I really did not have enough time at the end (I should, for instance, have had three or four more lectures in order to deal more completely with such matters as energy bands and the spatial dependence of amplitudes). Also, I had never presented the subject this way before, so the lack of feedback was particularly serious. I now believe the quantum mechanics should be given at a later time. Maybe I'll have a chance to do it again someday. Then I'll do it right.

The reason there are no lectures on how to solve problems is because there were recitation sections. Although I did put in three lectures in the first year on how to solve problems, they are not included here. Also there was a lecture on inertial

guidance which certainly belongs after the lecture on rotating systems, but which was, unfortunately, omitted. The fifth and sixth lectures are actually due to Matthew Sands, as I was out of town.

The question, of course, is how well this experiment has succeeded. My own point of view—which, however, does not seem to be shared by most of the people who worked with the students—is pessimistic. I don't think I did very well by the students. When I look at the way the majority of the students handled the problems on the examinations, I think that the system is a failure. Of course, my friends point out to me that there were one or two dozen students who—very surprisingly—understood almost everything in all of the lectures, and who were quite active in working with the material and worrying about the many points in an excited and interested way. These people have now, I believe, a first-rate background in physics—and they are, after all, the ones I was trying to get at. But then, “The power of instruction is seldom of much efficacy except in those happy dispositions where it is almost superfluous.” (Gibbons)

Still, I didn't want to leave any student completely behind, as perhaps I did. I think one way we could help the students more would be by putting more hard work into developing a set of problems which would elucidate some of the ideas in the lectures. Problems give a good opportunity to fill out the material of the lectures and make more realistic, more complete, and more settled in the mind the ideas that have been exposed.

I think, however, that there isn't any solution to this problem of education other than to realize that the best teaching can be done only when there is a direct individual relationship between a student and a good teacher—a situation in which the student discusses the ideas, thinks about the things, and talks about the things. It's impossible to learn very much by simply sitting in a lecture, or even by simply doing problems that are assigned. But in our modern times we have so many students to teach that we have to try to find some substitute for the ideal. Perhaps my lectures can make some contribution. Perhaps in some small place where there are individual teachers and students, they may get some inspiration or some ideas from the lectures. Perhaps they will have fun thinking them through—or going on to develop some of the ideas further.

RICHARD P. FEYNMAN

June, 1963

Foreword

For some forty years Richard P. Feynman focussed his curiosity on the mysterious workings of the physical world, and bent his intellect to searching out the order in its chaos. Now, he has given two years of his ability and his energy to his Lectures on Physics for beginning students. For them he has distilled the essence of his knowledge, and has created in terms they can hope to grasp a picture of the physicist's universe. To his lectures he has brought the brilliance and clarity of his thought, the originality and vitality of his approach, and the contagious enthusiasm of his delivery. It was a joy to behold.

The first year's lectures formed the basis for the first volume of this set of books. We have tried in this the second volume to make some kind of a record of a part of the second year's lectures—which were given to the sophomore class during the 1962–1963 academic year. The rest of the second year's lectures will make up Volume III.

Of the second year of lectures, the first two-thirds were devoted to a fairly complete treatment of the physics of electricity and magnetism. Its presentation was intended to serve a dual purpose. We hoped, first, to give the students a complete view of one of the great chapters of physics—from the early gropings of Franklin, through the great synthesis of Maxwell, on to the Lorentz electron theory of material properties, and ending with the still unsolved dilemmas of the electromagnetic self-energy. And we hoped, second, by introducing at the outset the calculus of vector fields, to give a solid introduction to the mathematics of field theories. To emphasize the general utility of the mathematical methods, related subjects from other parts of physics were sometimes analyzed together with their electric counterparts. We continually tried to drive home the generality of the mathematics. (“The same equations have the same solutions.”) And we emphasized this point by the kinds of exercises and examinations we gave with the course.

Following the electromagnetism there are two chapters each on elasticity and fluid flow. In the first chapter of each pair, the elementary and practical aspects are treated. The second chapter on each subject attempts to give an overview of the whole complex range of phenomena which the subject can lead to. These four chapters can well be omitted without serious loss, since they are not at all a necessary preparation for Volume III.

The last quarter, approximately, of the second year was dedicated to an introduction to quantum mechanics. This material has been put into the third volume.

In this record of the Feynman Lectures we wished to do more than provide a transcription of what was said. We hoped to make the written version as clear an exposition as possible of the ideas on which the original lectures were based. For some of the lectures this could be done by making only minor adjustments of the wording in the original transcript. For others of the lectures a major reworking and rearrangement of the material was required. Sometimes we felt we should add some new material to improve the clarity or balance of the presentation. Throughout the process we benefitted from the continual help and advice of Professor Feynman.

The translation of over 1,000,000 spoken words into a coherent text on a tight schedule is a formidable task, particularly when it is accompanied by the

other onerous burdens which come with the introduction of a new course—preparing for recitation sections, and meeting students, designing exercises and examinations, and grading them, and so on. Many hands—and heads—were involved. In some instances we have, I believe, been able to render a faithful image—or a tenderly retouched portrait—of the original Feynman. In other instances we have fallen far short of this ideal. Our successes are owed to all those who helped. The failures, we regret.

As explained in detail in the Foreword to Volume I, these lectures were but one aspect of a program initiated and supervised by the Physics Course Revision Committee (R. B. Leighton, Chairman, H. V. Neher, and M. Sands) at the California Institute of Technology, and supported financially by the Ford Foundation. In addition, the following people helped with one aspect or another of the preparation of textual material for this second volume: T. K. Caughey, M. L. Clayton, J. B. Curcio, J. B. Hartle, T. W. H. Harvey, M. H. Israel, W. J. Karzas, R. W. Kavanagh, R. B. Leighton, J. Mathews, M. S. Plesset, F. L. Warren, W. Whaling, C. H. Wilts, and B. Zimmerman. Others contributed indirectly through their work on the course: J. Blue, G. F. Chapline, M. J. Clauser, R. Dolen, H. H. Hill, and A. M. Title. Professor Gerry Neugebauer contributed in all aspects of our task with a diligence and devotion far beyond the dictates of duty.

The story of physics you find here would, however, not have been, except for the extraordinary ability and industry of Richard P. Feynman.

MATTHEW SANDS

March, 1964

Contents

CHAPTER 1. ELECTROMAGNETISM

- 1-1 Electrical forces 1-1
- 1-2 Electric and magnetic fields 1-3
- 1-3 Characteristics of vector fields 1-4
- 1-4 The laws of electromagnetism 1-5
- 1-5 What are the fields? 1-9
- 1-6 Electromagnetism in science and technology 1-10

CHAPTER 2. DIFFERENTIAL CALCULUS OF VECTOR FIELDS

- 2-1 Understanding physics 2-1
- 2-2 Scalar and vector fields— T and h 2-2
- 2-3 Derivatives of fields—the gradient 2-4
- 2-4 The operator ∇ 2-6
- 2-5 Operations with ∇ 2-7
- 2-6 The differential equation of heat flow 2-8
- 2-7 Second derivatives of vector fields 2-9
- 2-8 Pitfalls 2-11

CHAPTER 3. VECTOR INTEGRAL CALCULUS

- 3-1 Vector integrals; the line integral of $\nabla\psi$ 3-1
- 3-2 The flux of a vector field 3-2
- 3-3 The flux from a cube; Gauss' theorem 3-4
- 3-4 Heat conduction; the diffusion equation 3-6
- 3-5 The circulation of a vector field 3-8
- 3-6 The circulation around a square; Stokes' theorem 3-9
- 3-7 Curl-free and divergence-free fields 3-10
- 3-8 Summary 3-11

CHAPTER 4. ELECTROSTATICS

- 4-1 Statics 4-1
- 4-2 Coulomb's law; superposition 4-2
- 4-3 Electric potential 4-4
- 4-4 $E = -\nabla\phi$ 4-6
- 4-5 The flux of E 4-7
- 4-6 Gauss' law; divergence of E 4-9
- 4-7 Field of a sphere of charge 4-10
- 4-8 Field lines; equipotential surfaces 4-11

CHAPTER 5. APPLICATION OF GAUSS' LAW

- 5-1 Electrostatics is Gauss's law plus . . . 5-1
- 5-2 Equilibrium in an electrostatic field 5-1
- 5-3 Equilibrium with conductors 5-2
- 5-4 Stability of atoms 5-3
- 5-5 The field of a line charge 5-3
- 5-6 A sheet of charge; two sheets 5-4
- 5-7 A sphere of charge; a spherical shell 5-4
- 5-8 Is the field of a point charge exactly $1/r^2$? 5-5
- 5-9 The fields of a conductor 5-7
- 5-10 The field in a cavity of a conductor 5-8

CHAPTER 6. THE ELECTRIC FIELD IN VARIOUS CIRCUMSTANCES

- 6-1 Equations of the electrostatic potential 6-1
- 6-2 The electric dipole 6-2
- 6-3 Remarks on vector equations 6-4
- 6-4 The dipole potential as a gradient 6-4
- 6-5 The dipole approximation for an arbitrary distribution 6-6
- 6-6 The fields of charged conductors 6-8
- 6-7 The method of images 6-8
- 6-8 A point charge near a conducting plane 6-9
- 6-9 A point charge near a conducting sphere 6-10
- 6-10 Condensers; parallel plates 6-11
- 6-11 High-voltage breakdown 6-13
- 6-12 The field-emission microscope 6-14

CHAPTER 7. THE ELECTRIC FIELD IN VARIOUS CIRCUMSTANCES (Continued)

- 7-1 Methods for finding the electrostatic field 7-1
- 7-2 Two-dimensional fields; functions of the complex variable 7-2
- 7-3 Plasma oscillations 7-5
- 7-4 Colloidal particles in an electrolyte 7-8
- 7-5 The electrostatic field of a grid 7-10

CHAPTER 8. ELECTROSTATIC ENERGY

- 8-1 The electrostatic energy of charges. A uniform sphere 8-1
- 8-2 The energy of a condenser. Forces on charged conductors 8-2
- 8-3 The electrostatic energy of an ionic crystal 8-4
- 8-4 Electrostatic energy in nuclei 8-6
- 8-5 Energy in the electrostatic field 8-9
- 8-6 The energy of a point charge 8-12

CHAPTER 9. ELECTRICITY IN THE ATMOSPHERE

- 9-1 The electric potential gradient of the atmosphere 9-1
- 9-2 Electric currents in the atmosphere 9-2
- 9-3 Origin of the atmospheric currents 9-4
- 9-4 Thunderstorms 9-5
- 9-5 The mechanism of charge separation 9-7
- 9-6 Lightning 9-10

CHAPTER 10. DIELECTRICS

- 10-1 The dielectric constant 10-1
- 10-2 The polarization vector P 10-2
- 10-3 Polarization charges 10-3
- 10-4 The electrostatic equations with dielectrics 10-6
- 10-5 Fields and forces with dielectrics 10-7

CHAPTER 11. INSIDE DIELECTRICS

- 11-1 Molecular dipoles 11-1
- 11-2 Electronic polarization 11-1
- 11-3 Polar molecules; orientation polarization 11-3
- 11-4 Electric fields in cavities of a dielectric 11-5
- 11-5 The dielectric constant of liquids; the Clausius-Mossotti equation 11-6
- 11-6 Solid dielectrics 11-8
- 11-7 Ferroelectricity; BaTiO_3 11-8

CHAPTER 12. ELECTROSTATIC ANALOGS

- 12-1 The same equations have the same solutions 12-1
- 12-2 The flow of heat; a point source near an infinite plane boundary 12-2
- 12-3 The stretched membrane 12-5
- 12-4 The diffusion of neutrons; a uniform spherical source in a homogeneous medium 12-6
- 12-5 Irrotational fluid flow; the flow past a sphere 12-8
- 12-6 Illumination; the uniform lighting of a plane 12-10
- 12-7 The "underlying unity" of nature 12-12

CHAPTER 13. MAGNETOSTATICS

- 13-1 The magnetic field 13-1
- 13-2 Electric current; the conservation of charge 13-1
- 13-3 The magnetic force on a current 13-2
- 13-4 The magnetic field of steady currents; Ampere's law 13-3
- 13-5 The magnetic field of a straight wire and of a solenoid; atomic currents 13-5
- 13-6 The relativity of magnetic and electric fields 13-6
- 13-7 The transformation of currents and charges 13-11
- 13-8 Superposition; the right-hand rule 13-11

CHAPTER 14. THE MAGNETIC FIELD IN VARIOUS SITUATIONS

- 14-1 The vector potential 14-1
- 14-2 The vector potential of known currents 14-3
- 14-3 A straight wire 14-4
- 14-4 A long solenoid 14-5
- 14-5 The field of a small loop; the magnetic dipole 14-7
- 14-6 The vector potential of a circuit 14-8
- 14-7 The law of Biot and Savart 14-9

CHAPTER 15. THE VECTOR POTENTIAL

- 15-1 The forces on a current loop; energy of a dipole 15-1
- 15-2 Mechanical and electrical energies 15-3
- 15-3 The energy of steady currents 15-6
- 15-4 B versus A 15-7
- 15-5 The vector potential and quantum mechanics 15-8
- 15-6 What is true for statics is false for dynamics 15-14

CHAPTER 16. INDUCED CURRENTS

- 16-1 Motors and generators 16-1
- 16-2 Transformers and inductances 16-4
- 16-3 Forces on induced currents 16-5
- 16-4 Electrical technology 16-8

CHAPTER 17. THE LAWS OF INDUCTION

- 17-1 The physics of induction 17-1
- 17-2 Exceptions to the "flux rule" 17-2
- 17-3 Particle acceleration by an induced electric field; the betatron 17-3
- 17-4 A paradox 17-5
- 17-5 Alternating-current generator 17-6
- 17-6 Mutual inductance 17-9
- 17-7 Self-inductance 17-11
- 17-8 Inductance and magnetic energy 17-12

CHAPTER 18. THE MAXWELL EQUATIONS

- 18-1 Maxwell's equations 18-1
- 18-2 How the new term works 18-3
- 18-3 All of classical physics 18-5
- 18-4 A travelling field 18-5
- 18-5 The speed of light 18-8
- 18-6 Solving Maxwell's equations; the potentials and the wave equation 18-9

CHAPTER 19. THE PRINCIPLE OF LEAST ACTION

- A special lecture—almost verbatim 19-1
- A note added after the lecture 19-14

CHAPTER 20. SOLUTIONS OF MAXWELL'S EQUATIONS IN FREE SPACE

- 20-1 Waves in free space; plane waves 20-1
- 20-2 Three-dimensional waves 20-8
- 20-3 Scientific imagination 20-9
- 20-4 Spherical waves 20-12

CHAPTER 21. SOLUTIONS OF MAXWELL'S EQUATIONS WITH CURRENTS AND CHARGES

- 21-1 Light and electromagnetic waves 21-1
- 21-2 Spherical waves from a point source 21-2
- 21-3 The general solution of Maxwell's equations 21-4
- 21-4 The fields of an oscillating dipole 21-5
- 21-5 The potentials of a moving charge; the general solution of Liénard and Wiechert 21-9
- 21-6 The potentials for a charge moving with constant velocity; the Lorentz formula 21-12

CHAPTER 22. AC CIRCUITS

- 22-1 Impedances 22-1
- 22-2 Generators 22-5
- 22-3 Networks of ideal elements; Kirchhoff's rules 22-7
- 22-4 Equivalent circuits 22-10
- 22-5 Energy 22-11
- 22-6 A ladder network 22-12
- 22-7 Filters 22-14
- 22-8 Other circuit elements 22-16

CHAPTER 23. CAVITY RESONATORS

- 23-1 Real circuit elements 23-1
- 23-2 A capacitor at high frequencies 23-2
- 23-3 A resonant cavity 23-6
- 23-4 Cavity modes 23-9
- 23-5 Cavities and resonant circuits 23-10

CHAPTER 24. WAVEGUIDES

- 24-1 The transmission line 24-1
- 24-2 The rectangular waveguide 24-4
- 24-3 The cutoff frequency 24-6
- 24-4 The speed of the guided waves 24-7
- 24-5 Observing guided waves 24-7
- 24-6 Waveguide plumbing 24-8
- 24-7 Waveguide modes 24-10
- 24-8 Another way of looking at the guided waves 24-10

CHAPTER 25. ELECTRODYNAMICS IN RELATIVISTIC NOTATION

- 25-1 Four-vectors 25-1
- 25-2 The scalar product 25-3
- 25-3 The four-dimensional gradient 25-6
- 25-4 Electrodynamics in four-dimensional notation 25-8
- 25-5 The four-potential of a moving charge 25-9
- 25-6 The invariance of the equations of electrodynamics 25-10

CHAPTER 26. LORENTZ TRANSFORMATIONS OF THE FIELDS

- 26-1 The four-potential of a moving charge 26-1
- 26-2 The fields of a point charge with a constant velocity 26-2
- 26-3 Relativistic transformation of the fields 26-5
- 26-4 The equations of motion in relativistic notation 26-11

CHAPTER 27. FIELD ENERGY AND FIELD MOMENTUM

- 27-1 Local conservation 27-1
- 27-2 Energy conservation and electromagnetism 27-2
- 27-3 Energy density and energy flow in the electromagnetic field 27-3
- 27-4 The ambiguity of the field energy 27-6
- 27-5 Examples of energy flow 27-6
- 27-6 Field momentum 27-9

CHAPTER 28. ELECTROMAGNETIC MASS

- 28-1 The field energy of a point charge 28-1
- 28-2 The field momentum of a moving charge 28-2
- 28-3 Electromagnetic mass 28-3
- 28-4 The force of an electron on itself 28-4
- 28-5 Attempts to modify the Maxwell theory 28-6
- 28-6 The nuclear force field 28-12

CHAPTER 29. THE MOTION OF CHARGES IN ELECTRIC AND MAGNETIC FIELDS

- 29-1 Motion in a uniform electric or magnetic field 29-1
- 29-2 Momentum analysis 29-1
- 29-3 An electrostatic lens 29-2
- 29-4 A magnetic lens 29-3
- 29-5 The electron microscope 29-3
- 29-6 Accelerator guide fields 29-4
- 29-7 Alternating-gradient focusing 29-6
- 29-8 Motion in crossed electric and magnetic fields 29-8

CHAPTER 30. THE INTERNAL GEOMETRY OF CRYSTALS

- 30-1 The internal geometry of crystals 30-1
- 30-2 Chemical bonds in crystals 30-2
- 30-3 The growth of crystals 30-3
- 30-4 Crystal lattices 30-3
- 30-5 Symmetries in two dimensions 30-4
- 30-6 Symmetries in three dimensions 30-7
- 30-7 The strength of metals 30-8
- 30-8 Dislocations and crystal growth 30-9
- 30-9 The Bragg-Nye crystal model 30-10

CHAPTER 31. TENSORS

- 31-1 The tensor of polarizability 31-1
- 31-2 Transforming the tensor components 31-3
- 31-3 The energy ellipsoid 31-3
- 31-4 Other tensors; the tensor of inertia 31-6
- 31-5 The cross product 31-8
- 31-6 The tensor of stress 31-9
- 31-7 Tensors of higher rank 31-11
- 31-8 The four-tensor of electromagnetic momentum 31-12

CHAPTER 32. REFRACTIVE INDEX OF DENSE MATERIALS

- 32-1 Polarization of matter 32-1
- 32-2 Maxwell's equations in a dielectric 32-3
- 32-3 Waves in a dielectric 32-5
- 32-4 The complex index of refraction 32-8
- 32-5 The index of a mixture 32-8
- 32-6 Waves in metals 32-10
- 32-7 Low-frequency and high-frequency approximations; the skin depth and the plasma frequency 32-11

CHAPTER 33. REFLECTION FROM SURFACES

- 33-1 Reflection and refraction of light 33-1
- 33-2 Waves in dense materials 33-2
- 33-3 The boundary conditions 33-4
- 33-4 The reflected and transmitted waves 33-7
- 33-5 Reflection from metals 33-11
- 33-6 Total internal reflection 33-12

CHAPTER 34. THE MAGNETISM OF MATTER

- 34-1 Diamagnetism and paramagnetism 34-1
- 34-2 Magnetic moments and angular momentum 34-3
- 34-3 The precession of atomic magnets 34-4
- 34-4 Diamagnetism 34-5
- 34-5 Larmor's theorem 34-6
- 34-6 Classical physics gives neither diamagnetism nor paramagnetism 34-8
- 34-7 Angular momentum in quantum mechanics 34-8
- 34-8 The magnetic energy of atoms 34-11

CHAPTER 35. PARAMAGNETISM AND MAGNETIC RESONANCE

- 35-1 Quantized magnetic states 35-1
- 35-2 The Stern-Gerlach experiment 35-3
- 35-3 The Rabi molecular-beam method 35-4
- 35-4 The paramagnetism of bulk materials 35-6
- 35-5 Cooling by adiabatic demagnetization 35-9
- 35-6 Nuclear magnetic resonance 35-10

CHAPTER 36. FERROMAGNETISM

- 36-1 Magnetization currents 36-1
- 36-2 The field H 36-5
- 36-3 The magnetization curve 36-6
- 36-4 Iron-core inductances 36-8
- 36-5 Electromagnets 36-9
- 36-6 Spontaneous magnetization 36-11

CHAPTER 37. MAGNETIC MATERIALS

- 37-1 Understanding ferromagnetism 37-1
- 37-2 Thermodynamic properties 37-4
- 37-3 The hysteresis curve 37-5
- 37-4 Ferromagnetic materials 37-10
- 37-5 Extraordinary magnetic materials 37-11

CHAPTER 38. ELASTICITY

- 38-1 Hooke's law 38-1
- 38-2 Uniform strains 38-2
- 38-3 The torsion bar; shear waves 38-5
- 38-4 The bent beam 38-9
- 38-5 Buckling 38-11

CHAPTER 39. ELASTIC MATERIALS

- 39-1 The tensor of strain 39-1
- 39-2 The tensor of elasticity 39-4
- 39-3 The motions in an elastic body 39-6
- 39-4 Nonelastic behavior 39-8
- 39-5 Calculating the elastic constants 39-10

CHAPTER 40. THE FLOW OF DRY WATER

- 40-1 Hydrostatics 40-1
- 40-2 The equations of motion 40-2
- 40-3 Steady flow—Bernoulli's theorem 40-6
- 40-4 Circulation 40-9
- 40-5 Vortex lines 40-10

CHAPTER 41. THE FLOW OF WET WATER

- 41-1 Viscosity 41-1
- 41-2 Viscous flow 41-4
- 41-3 The Reynolds number 41-5
- 41-4 Flow past a circular cylinder 41-7
- 41-5 The limit of zero viscosity 41-9
- 41-6 Couette flow 41-10

INDEX

Electromagnetism

1-1 Electrical forces

Consider a force like gravitation which varies predominantly inversely as the square of the distance, but which is about a *billion-billion-billion-billion* times stronger. And with another difference. There are two kinds of "matter," which we can call positive and negative. Like kinds repel and unlike kinds attract—unlike gravity where there is only attraction. What would happen?

A bunch of positives would repel with an enormous force and spread out in all directions. A bunch of negatives would do the same. But an evenly mixed bunch of positives and negatives would do something completely different. The opposite pieces would be pulled together by the enormous attractions. The net result would be that the terrific forces would balance themselves out almost perfectly, by forming tight, fine mixtures of the positive and the negative, and between two separate bunches of such mixtures there would be practically no attraction or repulsion at all.

There is such a force: the electrical force. And all matter is a mixture of positive protons and negative electrons which are attracting and repelling with this great force. So perfect is the balance, however, that when you stand near someone else you don't feel any force at all. If there were even a little bit of unbalance you would know it. If you were standing at arm's length from someone and each of you had *one percent* more electrons than protons, the repelling force would be incredible. How great? Enough to lift the Empire State Building? No! To lift Mount Everest? No! The repulsion would be enough to lift a "weight" equal to that of the entire earth!

With such enormous forces so perfectly balanced in this intimate mixture, it is not hard to understand that matter, trying to keep its positive and negative charges in the finest balance, can have a great stiffness and strength. The Empire State Building, for example, swings only eight feet in the wind because the electrical forces hold every electron and proton more or less in its proper place. On the other hand, if we look at matter on a scale small enough that we see only a few atoms, any small piece will not, usually, have an equal number of positive and negative charges, and so there will be strong residual electrical forces. Even when there are equal numbers of both charges in two neighboring small pieces, there may still be large net electrical forces because the forces between individual charges vary inversely as the square of the distance. A net force can arise if a negative charge of one piece is closer to the positive than to the negative charges of the other piece. The attractive forces can then be larger than the repulsive ones and there can be a net attraction between two small pieces with no excess charges. The force that holds the atoms together, and the chemical forces that hold molecules together, are really electrical forces acting in regions where the balance of charge is not perfect, or where the distances are very small.

You know, of course, that atoms are made with positive protons in the nucleus and with electrons outside. You may ask: "If this electrical force is so terrific, why don't the protons and electrons just get on top of each other? If they want to be in an intimate mixture, why isn't it still more intimate?" The answer has to do with the quantum effects. If we try to confine our electrons in a region that is very close to the protons, then according to the uncertainty principle they must have some mean square momentum which is larger the more we try to confine them. It is this motion, required by the laws of quantum mechanics, that keeps the electrical attraction from bringing the charges any closer together.

1-1 Electrical forces

1-2 Electric and magnetic fields

1-3 Characteristics of vector fields

1-4 The laws of electromagnetism

1-5 What are the fields?

1-6 Electromagnetism in science and technology

Review: Chapter 12, Vol. I, Characteristics of Force

**Lower case Greek letters
and commonly used capitals**

α	alpha
β	beta
γ Γ	gamma
δ Δ	delta
ϵ	epsilon
ζ	zeta
η	eta
θ Θ	theta
ι	iota
κ	kappa
λ Λ	lambda
μ	mu
ν	nu
ξ Ξ	xi (ksi)
\omicron	omicron
π Π	pi
ρ	rho
σ Σ	sigma
τ	tau
υ Υ	upsilon
ϕ Φ	phi
χ	chi (khi)
ψ Ψ	psi
ω Ω	omega

There is another question: "What holds the nucleus together"? In a nucleus there are several protons, all of which are positive. Why don't they push themselves apart? It turns out that in nuclei there are, in addition to electrical forces, nonelectrical forces, called nuclear forces, which are greater than the electrical forces and which are able to hold the protons together in spite of the electrical repulsion. The nuclear forces, however, have a short range—their force falls off much more rapidly than $1/r^2$. And this has an important consequence. If a nucleus has too many protons in it, it gets too big, and it will not stay together. An example is uranium, with 92 protons. The nuclear forces act mainly between each proton (or neutron) and its nearest neighbor, while the electrical forces act over larger distances, giving a repulsion between each proton and all of the others in the nucleus. The more protons in a nucleus, the stronger is the electrical repulsion, until, as in the case of uranium, the balance is so delicate that the nucleus is almost ready to fly apart from the repulsive electrical force. If such a nucleus is just "tapped" lightly (as can be done by sending in a slow neutron), it breaks into two pieces, each with positive charge, and these pieces fly apart by electrical repulsion. The energy which is liberated is the energy of the atomic bomb. This energy is usually called "nuclear" energy, but it is really "electrical" energy released when electrical forces have overcome the attractive nuclear forces.

We may ask, finally, what holds a negatively charged electron together (since it has no nuclear forces). If an electron is all made of one kind of substance, each part should repel the other parts. Why, then, doesn't it fly apart? But does the electron have "parts"? Perhaps we should say that the electron is just a point and that electrical forces only act between *different* point charges, so that the electron does not act upon itself. Perhaps. All we can say is that the question of what holds the electron together has produced many difficulties in the attempts to form a complete theory of electromagnetism. The question has never been answered. We will entertain ourselves by discussing this subject some more in later chapters.

As we have seen, we should expect that it is a combination of electrical forces and quantum-mechanical effects that will determine the detailed structure of materials in bulk, and, therefore, their properties. Some materials are hard, some are soft. Some are electrical "conductors"—because their electrons are free to move about; others are "insulators"—because their electrons are held tightly to individual atoms. We shall consider later how some of these properties come about, but that is a very complicated subject, so we will begin by looking at the electrical forces only in simple situations. We begin by treating only the laws of electricity—including magnetism, which is really a part of the same subject.

We have said that the electrical force, like a gravitational force, decreases inversely as the square of the distance between charges. This relationship is called Coulomb's law. But it is not precisely true when charges are moving—the electrical forces depend also on the motions of the charges in a complicated way. One part of the force between moving charges we call the *magnetic* force. It is really one aspect of an electrical effect. That is why we call the subject "electromagnetism."

There is an important general principle that makes it possible to treat electromagnetic forces in a relatively simple way. We find, from experiment, that the force that acts on a particular charge—no matter how many other charges there are or how they are moving—depends only on the position of that particular charge, on the velocity of the charge, and on the amount of charge. We can write the force F on a charge q moving with a velocity v as

$$F = q(E + v \times B). \quad (1.1)$$

We call E the *electric field* and B the *magnetic field* at the location of the charge. The important thing is that the electrical forces from all the other charges in the universe can be summarized by giving just these two vectors. Their values will depend on *where* the charge is, and may change with *time*. Furthermore, if we replace that charge with another charge, the force on the new charge will be just in proportion to the amount of charge so long as all the rest of the charges in the

world do not change their positions or motions. (In real situations, of course, each charge produces forces on all other charges in the neighborhood and may cause these other charges to move, and so in some cases the fields *can* change if we replace our particular charge by another.)

We know from Vol. I how to find the motion of a particle if we know the force on it. Equation (1.1) can be combined with the equation of motion to give

$$\frac{d}{dt} \left[\frac{mv}{(1 - v^2/c^2)^{1/2}} \right] = \mathbf{F} = q(\mathbf{E} + \mathbf{v} \times \mathbf{B}). \quad (1.2)$$

So if \mathbf{E} and \mathbf{B} are given, we can find the motions. Now we need to know how the \mathbf{E} 's and \mathbf{B} 's are produced.

One of the most important simplifying principles about the way the fields are produced is this: Suppose a number of charges moving in some manner would produce a field \mathbf{E}_1 , and another set of charges would produce \mathbf{E}_2 . If both sets of charges are in place at the same time (keeping the same locations and motions they had when considered separately), then the field produced is just the sum

$$\mathbf{E} = \mathbf{E}_1 + \mathbf{E}_2. \quad (1.3)$$

This fact is called *the principle of superposition* of fields. It holds also for magnetic fields.

This principle means that if we know the law for the electric and magnetic fields produced by a *single* charge moving in an arbitrary way, then all the laws of electrodynamics are complete. If we want to know the force on charge A we need only calculate the \mathbf{E} and \mathbf{B} produced by each of the charges B, C, D , etc., and then add the \mathbf{E} 's and \mathbf{B} 's from all the charges to find the fields, and from them the forces acting on charge A . If it had only turned out that the field produced by a single charge was simple, this would be the neatest way to describe the laws of electrodynamics. We have already given a description of this law (Chapter 28, Vol. I) and it is, unfortunately, rather complicated.

It turns out that the form in which the laws of electrodynamics are simplest are not what you might expect. It is *not* simplest to give a formula for the force that one charge produces on another. It is true that when charges are standing still the Coulomb force law is simple, but when charges are moving about the relations are complicated by delays in time and by the effects of acceleration, among others. As a result, we do not wish to present electrodynamics only through the force laws between charges; we find it more convenient to consider another point of view—a point of view in which the laws of electrodynamics appear to be the most easily manageable.

1-2 Electric and magnetic fields

First, we must extend, somewhat, our ideas of the electric and magnetic vectors, \mathbf{E} and \mathbf{B} . We have defined them in terms of the forces that are felt by a charge. We wish now to speak of electric and magnetic fields *at a point* even when there is no charge present. We are saying, in effect, that since there are forces “acting on” the charge, there is still “something” there when the charge is removed. If a charge located at the point (x, y, z) at the time t feels the force \mathbf{F} given by Eq. (1.1) we associate the vectors \mathbf{E} and \mathbf{B} with *the point* in space (x, y, z) . We may think of $\mathbf{E}(x, y, z, t)$ and $\mathbf{B}(x, y, z, t)$ as giving the forces that *would be* experienced at the time t by a charge located at (x, y, z) , *with the condition* that placing the charge there *did not disturb* the positions or motions of all the other charges responsible for the fields.

Following this idea, we associate with *every* point (x, y, z) in space two vectors \mathbf{E} and \mathbf{B} , which may be changing with time. The electric and magnetic fields are, then, viewed as *vector functions* of x, y, z , and t . Since a vector is specified by its components, each of the fields $\mathbf{E}(x, y, z, t)$ and $\mathbf{B}(x, y, z, t)$ represent three mathematical functions of x, y, z , and t .

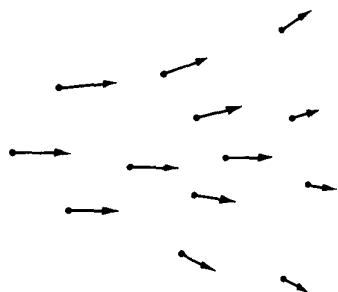


Fig. 1-1. A vector field may be represented by drawing a set of arrows whose magnitudes and directions indicate the values of the vector field at the points from which the arrows are drawn.

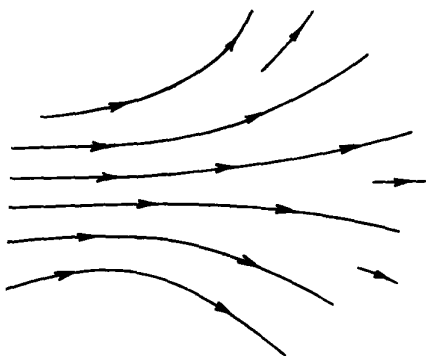


Fig. 1-2. A vector field can be represented by drawing lines which are tangent to the direction of the field vector at each point, and by drawing the density of lines proportional to the magnitude of the field vector.

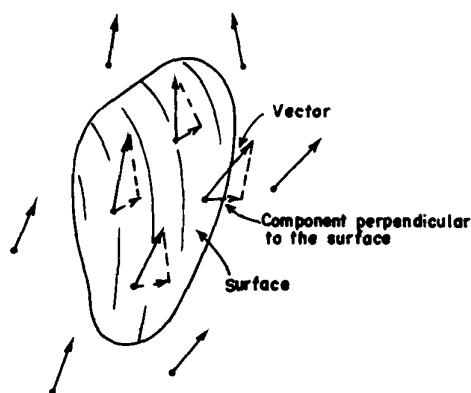


Fig. 1-3. The flux of a vector field through a surface is defined as the average value of the normal component of the vector times the area of the surface.

It is precisely because E (or B) can be specified at every point in space that it is called a "field." A "field" is any physical quantity which takes on different values at different points in space. Temperature, for example, is a field—in this case a scalar field, which we write as $T(x, y, z)$. The temperature could also vary in time, and we would say the temperature field is time-dependent, and write $T(x, y, z, t)$. Another example is the "velocity field" of a flowing liquid. We write $v(x, y, z, t)$ for the velocity of the liquid at each point in space at the time t . It is a vector field.

Returning to the electromagnetic fields—although they are produced by charges according to complicated formulas, they have the following important characteristic: the relationships between the values of the fields at *one point* and the values at a *nearby point* are very simple. With only a few such relationships in the form of differential equations we can describe the fields completely. It is in terms of such equations that the laws of electrodynamics are most simply written.

There have been various inventions to help the mind visualize the behavior of fields. The most correct is also the most abstract: we simply consider the fields as mathematical functions of position and time. We can also attempt to get a mental picture of the field by drawing vectors at many points in space, each of which gives the field strength and direction at that point. Such a representation is shown in Fig. 1-1. We can go further, however, and draw lines which are everywhere tangent to the vectors—which, so to speak, follow the arrows and keep track of the direction of the field. When we do this we lose track of the *lengths* of the vectors, but we can keep track of the strength of the field by drawing the lines far apart when the field is weak and close together when it is strong. We adopt the convention that the *number of lines per unit area* at right angles to the lines is proportional to the *field strength*. This is, of course, only an approximation, and it will require, in general, that new lines sometimes start up in order to keep the number up to the strength of the field. The field of Fig. 1-1 is represented by field lines in Fig. 1-2.

1-3 Characteristics of vector fields

There are two mathematically important properties of a vector field which we will use in our description of the laws of electricity from the field point of view. Suppose we imagine a closed surface of some kind and ask whether we are losing "something" from the inside; that is, does the field have a quality of "outflow"? For instance, for a velocity field we might ask whether the velocity is always outward on the surface or, more generally, whether more fluid flows out (per unit time) than comes in. We call the net amount of fluid going out through the surface per unit time the "flux of velocity" through the surface. The flow through an element of a surface is just equal to the component of the velocity perpendicular to the surface times the area of the surface. For an arbitrary closed surface, the *net outward flow*—or *flux*—is the average outward normal component of the velocity, times the area of the surface:

$$\text{Flux} = (\text{average normal component}) \cdot (\text{surface area}). \quad (1.4)$$

In the case of an electric field, we can mathematically define something analogous to an outflow, and we again call it the flux, but of course it is not the flow of any substance, because the electric field is not the velocity of anything. It turns out, however, that the mathematical quantity which is the average normal component of the field still has a useful significance. We speak, then, of the *electric flux*—also defined by Eq. (1.4). Finally, it is also useful to speak of the flux not only through a completely closed surface, but through any bounded surface. As before, the flux through such a surface is defined as the average normal component of a vector times the area of the surface. These ideas are illustrated in Fig. 1-3.

There is a second property of a vector field that has to do with a line, rather than a surface. Suppose again that we think of a velocity field that describes the flow of a liquid. We might ask this interesting question: Is the liquid circulating?

By that we mean: Is there a net rotational motion around some loop? Suppose that we instantaneously freeze the liquid everywhere except inside of a tube which is of uniform bore, and which goes in a loop that closes back on itself as in Fig. 1-4. Outside of the tube the liquid stops moving, but inside the tube it may keep on moving because of the momentum in the trapped liquid—that is, if there is more momentum heading one way around the tube than the other. We define a quantity called the *circulation* as the resulting speed of the liquid in the tube times its circumference. We can again extend our ideas and define the “circulation” for any vector field (even when there isn’t anything moving). For any vector field the *circulation around any imagined closed curve* is defined as the average tangential component of the vector (in a consistent sense) multiplied by the circumference of the loop (Fig. 1-5).

$$\text{Circulation} = (\text{average tangential component}) \cdot (\text{distance around}). \quad (1.5)$$

You will see that this definition does indeed give a number which is proportional to the circulation velocity in the quickly frozen tube described above.

With just these two ideas—flux and circulation—we can describe all the laws of electricity and magnetism at once. You may not understand the significance of the laws right away, but they will give you some idea of the way the physics of electromagnetism will be ultimately described.

1-4 The laws of electromagnetism

The first law of electromagnetism describes the flux of the electric field:

$$\text{The flux of } \mathbf{E} \text{ through any closed surface} = \frac{\text{the net charge inside}}{\epsilon_0}, \quad (1.6)$$

where ϵ_0 is a convenient constant. (The constant ϵ_0 is usually read as “epsilon-zero” or “epsilon-naught”.) If there are no charges inside the surface, even though there are charges nearby outside the surface, the *average* normal component of \mathbf{E} is zero, so there is no net flux through the surface. To show the power of this type of statement, we can show that Eq. (1.6) is the same as Coulomb’s law, provided only that we also add the idea that the field from a single charge is spherically symmetric. For a point charge, we draw a sphere around the charge. Then the average normal component is just the value of the magnitude of \mathbf{E} at any point, since the field must be directed radially and have the same strength for all points on the sphere. Our rule now says that the field at the surface of the sphere, times the area of the sphere—that is, the outgoing flux—is proportional to the charge inside. If we were to make the radius of the sphere bigger, the area would increase as the square of the radius. The average normal component of the electric field times that area must still be equal to the same charge inside, and so the field must decrease as the square of the distance—we get an “inverse square” field.

If we have an arbitrary curve in space and measure the circulation of the electric field around the curve, we will find that it is not, in general, zero (although it is for the Coulomb field). Rather, for electricity there is a second law that states: for any surface S (not closed) whose edge is the curve C ,

$$\text{Circulation of } \mathbf{E} \text{ around } C = \frac{d}{dt} (\text{flux of } \mathbf{B} \text{ through } S). \quad (1.7)$$

We can complete the laws of the electromagnetic field by writing two corresponding equations for the magnetic field \mathbf{B} .

$$\text{Flux of } \mathbf{B} \text{ through any closed surface} = 0. \quad (1.8)$$

For a surface S bounded by the curve C ,

$$c^2(\text{circulation of } \mathbf{B} \text{ around } C) = \frac{d}{dt} (\text{flux of } \mathbf{E} \text{ through } S) + \frac{\text{flux of electric current through } S}{\epsilon_0}. \quad (1.9)$$

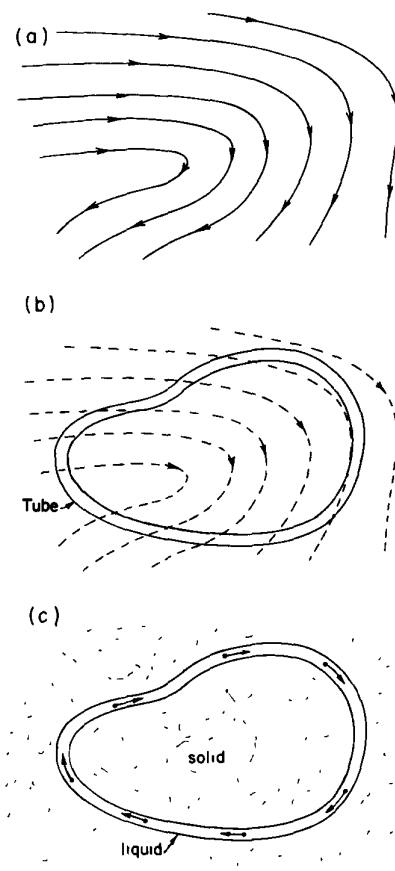


Fig. 1-4. (a) The velocity field in a liquid. Imagine a tube of uniform cross section that follows an arbitrary closed curve as in (b). If the liquid were suddenly frozen everywhere except inside the tube, the liquid in the tube would circulate as shown in (c).

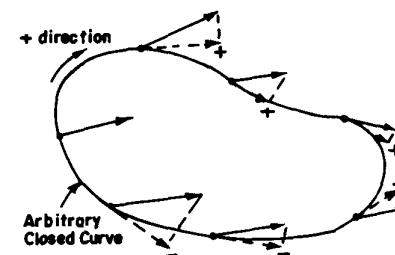


Fig. 1-5. The circulation of a vector field is the average tangential component of the vector (in a consistent sense) times the circumference of the loop.

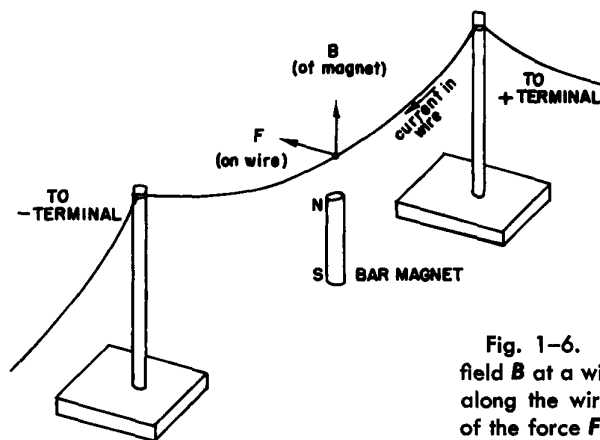


Fig. 1-6. A bar magnet gives a field B at a wire. When there is a current along the wire, the wire moves because of the force $F = qv \times B$.

The constant c^2 that appears in Eq. (1.9) is the square of the velocity of light. It appears because magnetism is in reality a relativistic effect of electricity. The constant ϵ_0 has been stuck in to make the units of electric current come out in a convenient way.

Equations (1.6) through (1.9), together with Eq. (1.1), are all the laws of electrodynamics*. As you remember, the laws of Newton were very simple to write down, but they had a lot of complicated consequences and it took us a long time to learn about them all. These laws are not nearly as simple to write down, which means that the consequences are going to be more elaborate and it will take us quite a lot of time to figure them all out.

We can illustrate some of the laws of electrodynamics by a series of small experiments which show qualitatively the interrelationships of electric and magnetic fields. You have experienced the first term of Eq. (1.1) when combing your hair, so we won't show that one. The second part of Eq. (1.1) can be demonstrated by passing a current through a wire which hangs above a bar magnet, as shown in Fig. 1-6. The wire will move when a current is turned on because of the force $F = qv \times B$. When a current exists, the charges inside the wire are moving, so they have a velocity v , and the magnetic field from the magnet exerts a force on them, which results in pushing the wire sideways.

When the wire is pushed to the left, we would expect that the magnet must feel a push to the right. (Otherwise we could put the whole thing on a wagon and have a propulsion system that didn't conserve momentum!) Although the force is too small to make movement of the bar magnet visible, a more sensitively supported magnet, like a compass needle, will show the movement.

How does the wire push on the magnet? The current in the wire produces a magnetic field of its own that exerts forces on the magnet. According to the last

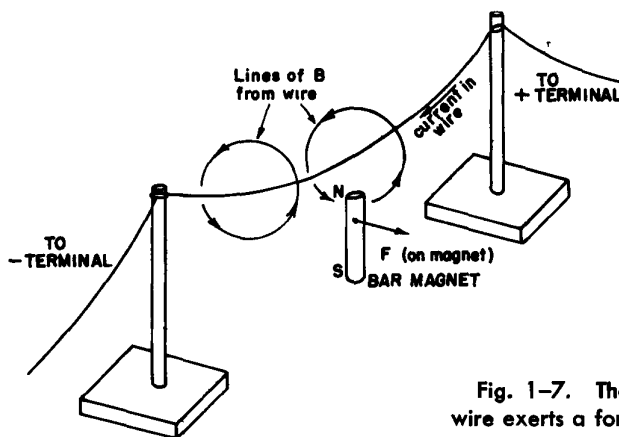


Fig. 1-7. The magnetic field of the wire exerts a force on the magnet.

* We need only to add a remark about some conventions for the *sign* of the circulation.

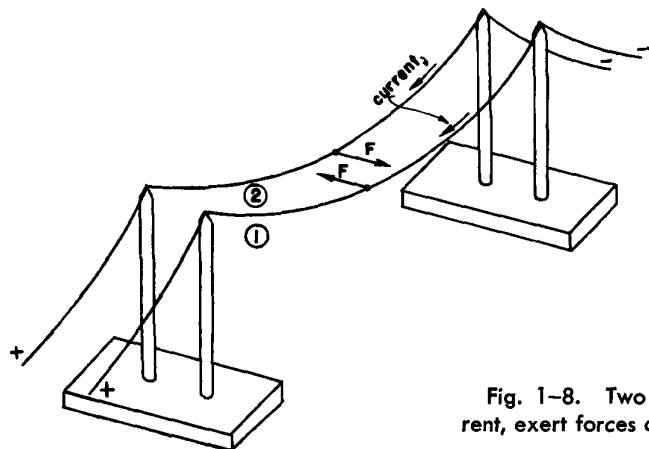


Fig. 1-8. Two wires, carrying current, exert forces on each other.

term in Eq. (1.9), a current must have a *circulation* of B —in this case, the lines of B are loops around the wire, as shown in Fig. 1-7. This B -field is responsible for the force on the magnet.

Equation (1.9) tells us that for a fixed current through the wire the circulation of B is the same for *any* curve that surrounds the wire. For curves—say circles—that are farther away from the wire, the circumference is larger, so the tangential component of B must decrease. You can see that we would, in fact, expect B to decrease linearly with the distance from a long straight wire.

Now, we have said that a current through a wire produces a magnetic field, and that when there is a magnetic field present there is a force on a wire carrying a current. Then we should also expect that if we make a magnetic field with a current in one wire, it should exert a force on another wire which also carries a current. This can be shown by using two hanging wires as shown in Fig. 1-8. When the currents are in the same direction, the two wires attract, but when the currents are opposite, they repel.

In short, electrical currents, as well as magnets, make magnetic fields. But wait, what is a magnet, anyway? If magnetic fields are produced by moving charges, is it not possible that the magnetic field from a piece of iron is really the result of currents? It appears to be so. We can replace the bar magnet of our experiment with a coil of wire, as shown in Fig. 1-9. When a current is passed through the coil—as well as through the straight wire above it—we observe a motion of the wire exactly as before, when we had a magnet instead of a coil. In other words, the current in the coil imitates a magnet. It appears, then, that a piece of iron acts as though it contains a perpetual circulating current. We can, in fact, understand magnets in terms of permanent currents in the atoms of the iron. The force on the magnet in Fig. 1-7 is due to the second term in Eq. (1.1).

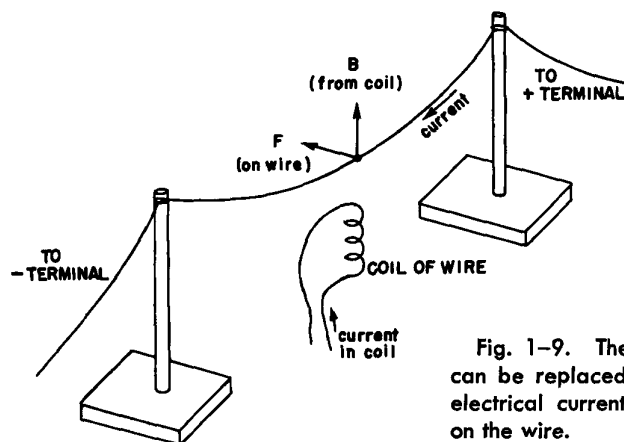
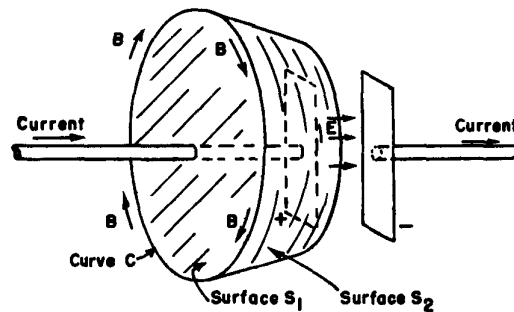


Fig. 1-9. The bar magnet of Fig. 1-6 can be replaced by a coil carrying an electrical current. A similar force acts on the wire.

Where do the currents come from? One possibility would be from the motion of the electrons in atomic orbits. Actually, that is not the case for iron, although it is for some materials. In addition to moving around in an atom, an electron also spins about on its own axis—something like the spin of the earth—and it is the current from this spin that gives the magnetic field in iron. (We say “something like the spin of the earth” because the question is so deep in quantum mechanics that the classical ideas do not really describe things too well.) In most substances, some electrons spin one way and some spin the other, so the magnetism cancels out, but in iron—for a mysterious reason which we will discuss later—many of the electrons are spinning with their axes lined up, and that is the source of the magnetism.

Since the fields of magnets are from currents, we do not have to add any extra term to Eqs. (1.8) or (1.9) to take care of magnets. We just take *all* currents, including the circulating currents of the spinning electrons, and then the law is right. You should also notice that Eq. (1.8) says that there are no magnetic “charges” analogous to the electrical charges appearing on the right side of Eq. (1.6). None has been found.

Fig. 1-10. The circulation of B around the curve C is given either by the current passing through the surface S_1 , or by the rate of change of the flux of E through the surface S_2 .



The first term on the right-hand side of Eq. (1.9) was discovered theoretically by Maxwell and is of great importance. It says that changing *electric* fields produce magnetic effects. In fact, without this term the equation would not make sense, because without it there could be no currents in circuits that are not complete loops. But such currents do exist, as we can see in the following example. Imagine a capacitor made of two flat plates. It is being charged by a current that flows toward one plate and away from the other, as shown in Fig. 1-10. We draw a curve C around one of the wires and fill it in with a surface which crosses the wire, as shown by the surface S_1 in the figure. According to Eq. (1.9), the circulation of B around C is given by the current in the wire (times c^2). But what if we fill in the curve with a *different* surface S_2 , which is shaped like a bowl and passes between the plates of the capacitor, staying always away from the wire? There is certainly no current through this surface. But, surely, just changing the location of an imaginary surface is not going to change a real magnetic field! The circulation of B must be what it was before. The first term on the right-hand side of Eq. (1.9) does, indeed, combine with the second term to give the same result for the two surfaces S_1 and S_2 . For S_2 the circulation of B is given in terms of the rate of change of the flux of E between the plates of the capacitor. And it works out that the changing E is related to the current in just the way required for Eq. (1.9) to be correct. Maxwell saw that it was needed, and he was the first to write the complete equation.

With the setup shown in Fig. 1-6 we can demonstrate another of the laws of electromagnetism. We disconnect the ends of the hanging wire from the battery and connect them to a galvanometer which tells us when there is a current through the wire. When we *push* the wire sideways through the magnetic field of the magnet, we observe a current. Such an effect is again just another consequence of Eq. (1.1)—the electrons in the wire feel the force $F = qv \times B$. The electrons have a sideways velocity because they move with the wire. This v with a vertical B from the magnet results in a force on the electrons directed *along* the wire, which starts the electrons moving toward the galvanometer.

Suppose, however, that we leave the wire alone and move the magnet. We guess from relativity that it should make no difference, and indeed, we observe a similar current in the galvanometer. How does the magnetic field produce forces on charges at rest? According to Eq. (1.1) there must be an electric field. A moving magnet must make an electric field. How that happens is said quantitatively by Eq. (1.7). This equation describes many phenomena of great practical interest, such as those that occur in electric generators and transformers.

The most remarkable consequence of our equations is that the combination of Eq. (1.7) and Eq. (1.9) contains the explanation of the radiation of electromagnetic effects over large distances. The reason is roughly something like this: suppose that somewhere we have a magnetic field which is increasing because, say, a current is turned on suddenly in a wire. Then by Eq. (1.7) there must be a circulation of an electric field. As the electric field builds up to produce its circulation, then according to Eq. (1.9) a magnetic circulation will be generated. But the building up of *this* magnetic field will produce a new circulation of the electric field, and so on. In this way fields work their way through space without the need of charges or currents except at their source. That is the way we *see* each other! It is all in the equations of the electromagnetic fields.

1-5 What are the fields?

We now make a few remarks on our way of looking at this subject. You may be saying: "All this business of fluxes and circulations is pretty abstract. There are electric fields at every point in space; then there are these 'laws.' But what is *actually* happening? Why can't you explain it, for instance, by whatever it *is* that goes between the charges." Well, it depends on your prejudices. Many physicists used to say that direct action with nothing in between was inconceivable. (How could they find an idea inconceivable when it had already been conceived?) They would say: "Look, the only forces we know are the direct action of one piece of matter on another. It is impossible that there can be a force with nothing to transmit it." But what really happens when we study the "direct action" of one piece of matter right against another? We discover that it is not one piece right against the other; they are slightly separated, and there are electrical forces acting on a tiny scale. Thus we find that we are going to explain so-called direct-contact action in terms of the picture for electrical forces. It is certainly not sensible to try to insist that an electrical force has to look like the old, familiar, muscular push or pull, when it will turn out that the muscular pushes and pulls are going to be interpreted as electrical forces! The only sensible question is what is the *most convenient* way to look at electrical effects. Some people prefer to represent them as the interaction at a distance of charges, and to use a complicated law. Others love the field lines. They draw field lines all the time, and feel that writing E 's and B 's is too abstract. The field lines, however, are only a crude way of describing a field, and it is very difficult to give the correct, quantitative laws directly in terms of field lines. Also, the ideas of the field lines do not contain the deepest principle of electrodynamics, which is the superposition principle. Even though we know how the field lines look for one set of charges and what the field lines look like for another set of charges, we don't get any idea about what the field line patterns will look like when both sets are present together. From the mathematical standpoint, on the other hand, superposition is easy—we simply add the two vectors. The field lines have some advantage in giving a vivid picture, but they also have some disadvantages. The direct interaction way of thinking has great advantages when thinking of electrical charges at rest, but has great disadvantages when dealing with charges in rapid motion.

The best way is to use the abstract field idea. That it is abstract is unfortunate, but necessary. The attempts to try to represent the electric field as the motion of some kind of gear wheels, or in terms of lines, or of stresses in some kind of material have used up more effort of physicists than it would have taken simply to get the right answers about electrodynamics. It is interesting that the correct equations for the behavior of light in crystals were worked out by McCullough in 1843. But

people said to him: "Yes, but there is no real material whose mechanical properties could possibly satisfy those equations, and since light is an oscillation that must vibrate in *something*, we cannot believe this abstract equation business." If people had been more open-minded, they might have believed in the right equations for the behavior of light a lot earlier than they did.

In the case of the magnetic field we can make the following point: Suppose that you finally succeeded in making up a picture of the magnetic field in terms of some kind of lines or of gear wheels running through space. Then you try to explain what happens to two charges moving in space, both at the same speed and parallel to each other. Because they are moving, they will behave like two currents and will have a magnetic field associated with them (like the currents in the wires of Fig. 1-8). An observer who was riding along with the two charges, however, would see both charges as stationary, and would say that there is *no* magnetic field. The "gear wheels" or "lines" disappear when you ride along with the object! All we have done is to invent a *new* problem. How can the gear wheels disappear?! The people who draw field lines are in a similar difficulty. Not only is it not possible to say whether the field lines move or do not move with charges—they may disappear completely in certain coordinate frames

What we are saying, then, is that magnetism is really a relativistic effect. In the case of the two charges we just considered, travelling parallel to each other, we would expect to have to make relativistic corrections to their motion, with terms of order v^2/c^2 . These corrections must correspond to the magnetic force. But what about the force between the two wires in our experiment (Fig. 1-8). There the magnetic force is the *whole* force. It didn't look like a "relativistic correction." Also, if we estimate the velocities of the electrons in the wire (you can do this yourself), we find that their average speed along the wire is about 0.01 centimeter per second. So v^2/c^2 is about 10^{-25} . Surely a negligible "correction." But no! Although the magnetic force is, in this case, 10^{-25} of the "normal" electrical force between the moving electrons, remember that the "normal" electrical forces have disappeared because of the almost perfect balancing out—because the wires have the same number of protons as electrons. The balance is much more precise than one part in 10^{25} , and the small relativistic term which we call the magnetic force is the only term left. It becomes the dominant term.

It is the near-perfect cancellation of electrical effects which allowed relativity effects (that is, magnetism) to be studied and the correct equations—to order v^2/c^2 —to be discovered, even though physicists didn't *know* that's what was happening. And that is why, when relativity was discovered, the electromagnetic laws didn't need to be changed. They—unlike mechanics—were already correct to a precision of v^2/c^2 .

1-6 Electromagnetism in science and technology

Let us end this chapter by pointing out that among the many phenomena studied by the Greeks there were two very strange ones: that if you rubbed a piece of amber you could lift up little pieces of papyrus, and that there was a strange rock from the island of Magnesia which attracted iron. It is amazing to think that these were the only phenomena known to the Greeks in which the effects of electricity or magnetism were apparent. The reason that these were the only phenomena that appeared is due primarily to the fantastic precision of the balancing of charges that we mentioned earlier. Study by scientists who came after the Greeks uncovered one new phenomena after another that were really some aspect of these amber and/or lodestone effects. Now we realize that the phenomena of chemical interaction and, ultimately, of life itself are to be understood in terms of electromagnetism.

At the same time that an understanding of the subject of electromagnetism was being developed, technical possibilities that defied the imagination of the people that came before were appearing: it became possible to signal by telegraph over long distances, and to talk to another person miles away without any connections between, and to run huge power systems—a great water wheel, connected by

filaments over hundreds of miles to another engine that turns in response to the master wheel—many thousands of branching filaments—ten thousand engines in ten thousand places running the machines of industries and homes—all turning because of the knowledge of the laws of electromagnetism.

Today we are applying even more subtle effects. The electrical forces, enormous as they are, can also be very tiny, and we can control them and use them in very many ways. So delicate are our instruments that we can tell what a man is doing by the way he affects the electrons in a thin metal rod hundreds of miles away. All we need to do is to use the rod as an antenna for a television receiver!

From a long view of the history of mankind—seen from, say, ten thousand years from now—there can be little doubt that the most significant event of the 19th century will be judged as Maxwell's discovery of the laws of electrodynamics. The American Civil War will pale into provincial insignificance in comparison with this important scientific event of the same decade.

Differential Calculus of Vector Fields

2-1 Understanding physics

The physicist needs a facility in looking at problems from several points of view. The exact analysis of real physical problems is usually quite complicated, and any particular physical situation may be too complicated to analyze directly by solving the differential equation. But one can still get a very good idea of the behavior of a system if one has some feel for the character of the solution in different circumstances. Ideas such as the field lines, capacitance, resistance, and inductance are, for such purposes, very useful. So we will spend much of our time analyzing them. In this way we will get a feel as to what should happen in different electromagnetic situations. On the other hand, none of the heuristic models, such as field lines, is really adequate and accurate for all situations. There is only one precise way of presenting the laws, and that is by means of differential equations. They have the advantage of being fundamental and, so far as we know, precise. If you have learned the differential equations you can always go back to them. There is nothing to unlearn.

It will take you some time to understand what should happen in different circumstances. You will have to solve the equations. Each time you solve the equations, you will learn something about the character of the solutions. To keep these solutions in mind, it will be useful also to study their meaning in terms of field lines and of other concepts. This is the way you will really “understand” the equations. That is the difference between mathematics and physics. Mathematicians, or people who have very mathematical minds, are often led astray when “studying” physics because they lose sight of the physics. They say: “Look, these differential equations—the Maxwell equations—are all there is to electrodynamics; it is admitted by the physicists that there is nothing which is not contained in the equations. The equations are complicated, but after all they are only mathematical equations and if I understand them mathematically inside out, I will understand the physics inside out.” Only it doesn’t work that way. Mathematicians who study physics with that point of view—and there have been many of them—usually make little contribution to physics and, in fact, little to mathematics. They fail because the actual physical situations in the real world are so complicated that it is necessary to have a much broader understanding of the equations.

What it means really to understand an equation—that is, in more than a strictly mathematical sense—was described by Dirac. He said: “I understand what an equation means if I have a way of figuring out the characteristics of its solution without actually solving it.” So if we have a way of knowing what should happen in given circumstances without actually solving the equations, then we “understand” the equations, as applied to these circumstances. A physical understanding is a completely unmathematical, imprecise, and inexact thing, but absolutely necessary for a physicist.

Ordinarily, a course like this is given by developing gradually the physical ideas—by starting with simple situations and going on to more and more complicated situations. This requires that you continuously forget things you previously learned—things that are true in certain situations, but which are not true in general. For example, the “law” that the electrical force depends on the square of the distance is not *always* true. We prefer the opposite approach. We prefer to take first the *complete* laws, and then to step back and apply them to simple situations, developing the physical ideas as we go along. And that is what we are going to do.

2-1 Understanding physics

2-2 Scalar and vector fields— T and h

2-3 Derivatives of fields—the gradient

2-4 The operator ∇

2-5 Operations with ∇

2-6 The differential equation of heat flow

2-7 Second derivatives of vector fields

2-8 Pitfalls

Review: Chapter 11, Vol. I, Vectors

Our approach is completely opposite to the historical approach in which one develops the subject in terms of the experiments by which the information was obtained. But the subject of physics has been developed over the past 200 years by some very ingenious people, and as we have only a limited time to acquire our knowledge, we cannot possibly cover everything they did. Unfortunately one of the things that we shall have a tendency to lose in these lectures is the historical, experimental development. It is hoped that in the laboratory some of this lack can be corrected. You can also fill in what we must leave out by reading the Encyclopedia Britannica, which has excellent historical articles on electricity and on other parts of physics. You will also find historical information in many textbooks on electricity and magnetism.

2-2 Scalar and vector fields— T and h

We begin now with the abstract, mathematical view of the theory of electricity and magnetism. The ultimate idea is to explain the meaning of the laws given in Chapter 1. But to do this we must first explain a new and peculiar notation that we want to use. So let us forget electromagnetism for the moment and discuss the mathematics of vector fields. It is of very great importance, not only for electromagnetism, but for all kinds of physical circumstances. Just as ordinary differential and integral calculus is so important to all branches of physics, so also is the differential calculus of vectors. We turn to that subject.

Listed below are a few facts from the algebra of vectors. It is assumed that you already know them.

$$\mathbf{A} \cdot \mathbf{B} = \text{scalar} = A_x B_x + A_y B_y + A_z B_z \quad (2.1)$$

$$\mathbf{A} \times \mathbf{B} = \text{vector} \quad (2.2)$$

$$(\mathbf{A} \times \mathbf{B})_z = A_x B_y - A_y B_x$$

$$(\mathbf{A} \times \mathbf{B})_x = A_y B_z - A_z B_y$$

$$(\mathbf{A} \times \mathbf{B})_y = A_z B_x - A_x B_z$$

$$\mathbf{A} \times \mathbf{A} = 0 \quad (2.3)$$

$$\mathbf{A} \cdot (\mathbf{A} \times \mathbf{B}) = 0 \quad (2.4)$$

$$\mathbf{A} \cdot (\mathbf{B} \times \mathbf{C}) = (\mathbf{A} \times \mathbf{B}) \cdot \mathbf{C} \quad (2.5)$$

$$\mathbf{A} \times (\mathbf{B} \times \mathbf{C}) = \mathbf{B}(\mathbf{A} \cdot \mathbf{C}) - \mathbf{C}(\mathbf{A} \cdot \mathbf{B}) \quad (2.6)$$

Also we will want to use the two following equalities from the calculus:

$$\Delta f(x, y, z) = \frac{\partial f}{\partial x} \Delta x + \frac{\partial f}{\partial y} \Delta y + \frac{\partial f}{\partial z} \Delta z, \quad (2.7)$$

$$\frac{\partial^2 f}{\partial x \partial y} = \frac{\partial^2 f}{\partial y \partial x}. \quad (2.8)$$

The first equation (2.7) is, of course, true only in the limit that Δx , Δy , and Δz go toward zero.

The simplest possible physical field is a scalar field. By a field, you remember, we mean a quantity which depends upon position in space. By a *scalar field* we merely mean a field which is characterized at each point by a single number—a scalar. Of course the number may change in time, but we need not worry about that for the moment. We will talk about what the field looks like at a given instant. As an example of a scalar field, consider a solid block of material which has been heated at some places and cooled at others, so that the temperature of the body varies from point to point in a complicated way. Then the temperature will be a function of x , y , and z , the position in space measured in a rectangular coordinate system. Temperature is a scalar field.

Writing vectors by hand.

Some people use

\vec{E} or \vec{E} or just \bar{E} .

Others prefer

\underline{E} .

We like the following way:

A B C D E F G

H I J K L M N

O P Q R S T U

V W X Y Z

Small letters are harder:

a b c d e f g

h i j k l m n

o p q r s t u

v w x y z

You can invent your own.

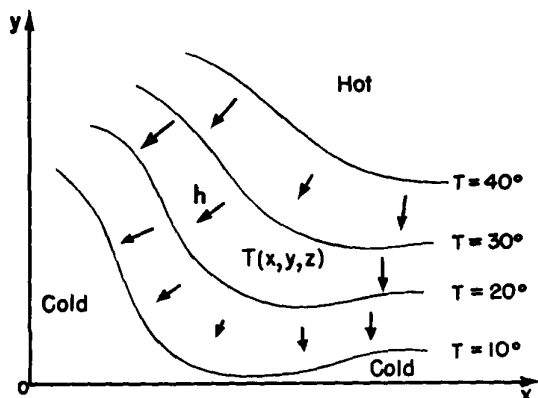


Fig. 2-1. Temperature T is an example of a scalar field. With each point (x, y, z) in space there is associated a number $T(x, y, z)$. All points on the surface marked $T = 20^\circ$ (shown as a curve at $z = 0$) are at the same temperature. The arrows are samples of the heat flow vector h .

One way of thinking about scalar fields is to imagine “contours” which are imaginary surfaces drawn through all points for which the field has the same value, just as contour lines on a map connect points with the same height. For a temperature field the contours are called “isothermal surfaces” or isotherms. Figure 2-1 illustrates a temperature field and shows the dependence of T on x and y when $z = 0$. Several isotherms are drawn.

There are also vector fields. The idea is very simple. A vector is given for each point in space. The vector varies from point to point. As an example, consider a rotating body. The velocity of the material of the body at any point is a vector which is a function of position (Fig. 2-2). As a second example, consider the flow of heat in a block of material. If the temperature in the block is high at one place and low at another, there will be a flow of heat from the hotter places to the colder. The heat will be flowing in different directions in different parts of the block. The heat flow is a directional quantity which we call h . Its magnitude is a measure of how much heat is flowing. Examples of the heat flow vector are also shown in Fig. 2-1.

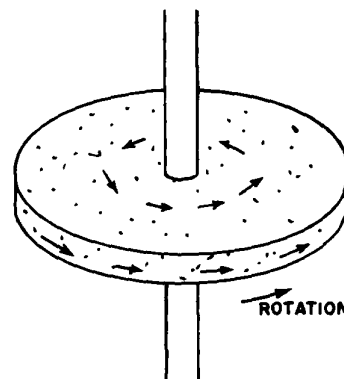


Fig. 2-2. The velocity of the atoms in a rotating object is an example of a vector field.

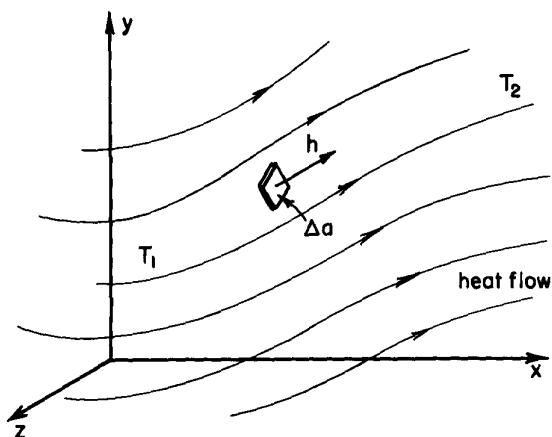


Fig. 2-3. Heat flow is a vector field. The vector h points along the direction of the flow. Its magnitude is the energy transported per unit time across a surface element oriented perpendicular to the flow, divided by the area of the surface element.

Let's make a more precise definition of h : The magnitude of the vector heat flow at a point is the amount of thermal energy that passes, per unit time and per unit area, through an infinitesimal surface element, at right angles to the direction of flow. The vector points in the direction of flow (see Fig. 2-3). In symbols: If ΔJ is the thermal energy that passes per unit time through the surface element Δa , then

$$h = \frac{\Delta J}{\Delta a} e_f, \quad (2.9)$$

where e_f is a unit vector in the direction of flow.

The vector h can be defined in another way—in terms of its components. We ask how much heat flows through a small surface at any angle with respect to the flow. In Fig. 2-4 we show a small surface Δa_2 inclined with respect to Δa_1 , which is perpendicular to the flow. The unit vector n is normal to the surface Δa_2 . The

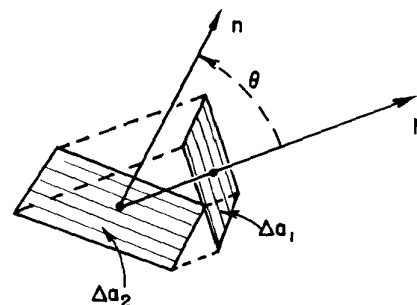


Fig. 2-4. The heat flow through Δa_2 is the same as through Δa_1 .

angle θ between \mathbf{n} and \mathbf{h} is the same as the angle between the surfaces (since \mathbf{h} is normal to Δa_1). Now what is the heat flow *per unit area* through Δa_2 ? The flow through Δa_2 is the same as through Δa_1 ; only the areas are different. In fact, $\Delta a_1 = \Delta a_2 \cos \theta$. The heat flow through Δa_2 is

$$\frac{\Delta J}{\Delta a_2} = \frac{\Delta J}{\Delta a_1} \cos \theta = \mathbf{h} \cdot \mathbf{n}. \quad (2.10)$$

We interpret this equation: the heat flow (per unit time and per unit area) through *any* surface element whose unit normal is \mathbf{n} , is given by $\mathbf{h} \cdot \mathbf{n}$. Equally, we could say: the component of the heat flow perpendicular to the surface element Δa_2 is $\mathbf{h} \cdot \mathbf{n}$. We can, if we wish, consider that these statements *define* \mathbf{h} . We will be applying the same ideas to other vector fields.

2-3 Derivatives of fields—the gradient

When fields vary in time, we can describe the variation by giving their derivatives with respect to t . We want to describe the variations with position in a similar way, because we are interested in the relationship between, say, the temperature in one place and the temperature at a nearby place. How shall we take the derivative of the temperature with respect to position? Do we differentiate the temperature with respect to x ? Or with respect to y , or z ?

Useful physical laws do not depend upon the orientation of the coordinate system. They should, therefore, be written in a form in which either both sides are scalars or both sides are vectors. What is the derivative of a scalar field, say $\partial T / \partial x$? Is it a scalar, or a vector, or what? It is neither a scalar nor a vector, as you can easily appreciate, because if we took a different x -axis, $\partial T / \partial x$ would certainly be different. But notice: We have three possible derivatives: $\partial T / \partial x$, $\partial T / \partial y$, and $\partial T / \partial z$. Since there are three kinds of derivatives and we know that it takes three numbers to form a vector, perhaps these three derivatives are the components of a vector:

$$\left(\frac{\partial T}{\partial x}, \frac{\partial T}{\partial y}, \frac{\partial T}{\partial z} \right) \stackrel{?}{=} \text{a vector}. \quad (2.11)$$

Of course it is not generally true that *any* three numbers form a vector. It is true only if, when we rotate the coordinate system, the components of the vector transform among themselves in the correct way. So it is necessary to analyze how these derivatives are changed by a rotation of the coordinate system. We shall show that (2.11) is indeed a vector. The derivatives do transform in the correct way when the coordinate system is rotated.

We can see this in several ways. One way is to ask a question whose answer is independent of the coordinate system, and try to express the answer in an “invariant” form. For instance, if $S = \mathbf{A} \cdot \mathbf{B}$, and if \mathbf{A} and \mathbf{B} are vectors, we know—because we proved it in Chapter 11 of Vol. I—that S is a scalar. We *know* that S is a scalar without investigating whether it changes with changes in coordinate systems. It *can't*, because it's a dot product of two vectors. Similarly, if we *know* that \mathbf{A} is a vector, and we have three numbers B_1 , B_2 , and B_3 , and we find out that

$$A_x B_1 + A_y B_2 + A_z B_3 = S, \quad (2.12)$$

where S is the same for any coordinate system, then it *must* be that the three numbers B_1 , B_2 , B_3 are the components B_x , B_y , B_z of some vector \mathbf{B} .

Now let's think of the temperature field. Suppose we take two points P_1 and P_2 , separated by the small interval ΔR . The temperature at P_1 is T_1 and at P_2 is T_2 , and the difference $\Delta T = T_2 - T_1$. The temperatures at these real, physical points certainly do not depend on what axis we choose for measuring the coordinates. In particular, ΔT is a number independent of the coordinate system. It is a scalar.

If we choose some convenient set of axes, we could write $T_1 = T(x, y, z)$ and $T_2 = T(x + \Delta x, y + \Delta y, z + \Delta z)$, where Δx , Δy , and Δz are the components of the vector $\Delta \mathbf{R}$ (Fig. 2-5). Remembering Eq. (2.7), we can write

$$\Delta T = \frac{\partial T}{\partial x} \Delta x + \frac{\partial T}{\partial y} \Delta y + \frac{\partial T}{\partial z} \Delta z. \quad (2.13)$$

The left side of Eq. (2.13) is a scalar. The right side is the sum of three products with Δx , Δy , and Δz , which are the components of a vector. It follows that the three numbers

$$\frac{\partial T}{\partial x}, \frac{\partial T}{\partial y}, \frac{\partial T}{\partial z}$$

are also the x -, y -, and z -components of a vector. We write this new vector with the symbol ∇T . The symbol ∇ (called “del”) is an upside-down Δ , and is supposed to remind us of differentiation. People read ∇T in various ways: “del- T ,” or “gradient of T ,” or “grad T ,”

$$\text{grad } T = \nabla T = \left(\frac{\partial T}{\partial x}, \frac{\partial T}{\partial y}, \frac{\partial T}{\partial z} \right)^*. \quad (2.14)$$

Using this notation, we can rewrite Eq. (2.13) in the more compact form

$$\Delta T = \nabla T \cdot \Delta \mathbf{R}. \quad (2.15)$$

In words, this equation says that the difference in temperature between two nearby points is the dot product of the gradient of T and the vector displacement between the points. The form of Eq. (2.15) also illustrates clearly our proof above that ∇T is indeed a vector.

Perhaps you are still not convinced? Let’s prove it in a different way. (Although if you look carefully, you may be able to see that it’s really the same proof in a longer-winded form!) We shall show that the components of ∇T transform in just the same way that components of \mathbf{R} do. If they do, ∇T is a vector according to our original definition of a vector in Chapter 11 of Vol. I. We take a new coordinate system x' , y' , z' , and in this new system we calculate $\partial T / \partial x'$, $\partial T / \partial y'$, and $\partial T / \partial z'$. To make things a little simpler, we let $z = z'$, so that we can forget about the z -coordinate. (You can check out the more general case for yourself.)

We take an $x'y'$ -system rotated an angle θ with respect to the xy -system, as in Fig. 2-6(a). For a point (x, y) the coordinates in the prime system are

$$x' = x \cos \theta + y \sin \theta, \quad (2.16)$$

$$y' = -x \sin \theta + y \cos \theta. \quad (2.17)$$

Or, solving for x and y ,

$$x = x' \cos \theta - y' \sin \theta, \quad (2.18)$$

$$y = x' \sin \theta + y' \cos \theta. \quad (2.19)$$

If any pair of numbers transforms with these equations in the same way that x and y do, they are the components of a vector.

Now let’s look at the difference in temperature between the two nearby points P_1 and P_2 , chosen as in Fig. 2-6(b). If we calculate with the x - and y -coordinates, we would write

$$\Delta T = \frac{\partial T}{\partial x} \Delta x \quad (2.20)$$

—since Δy is zero.

* In our notation, the expression (a, b, c) represents a vector with components a , b , and c . If you like to use the unit vectors i, j , and k , you may write

$$\nabla T = i \frac{\partial T}{\partial x} + j \frac{\partial T}{\partial y} + k \frac{\partial T}{\partial z}.$$

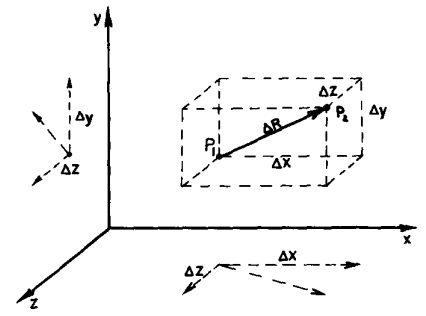


Fig. 2-5. The vector $\Delta \mathbf{R}$, whose components are Δx , Δy , and Δz .

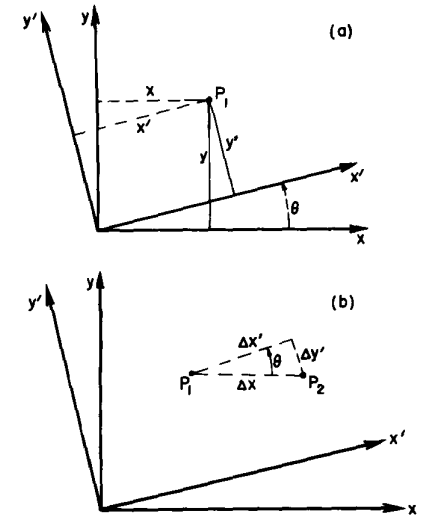


Fig. 2-6. (a) Transformation to a rotated coordinate system. (b) Special case of an interval $\Delta \mathbf{R}$ parallel to the x -axis.

If we choose some convenient set of axes, we could write $T_1 = T(x, y, z)$ and $T_2 = T(x + \Delta x, y + \Delta y, z + \Delta z)$, where Δx , Δy , and Δz are the components of the vector $\Delta \mathbf{R}$ (Fig. 2-5). Remembering Eq. (2.7), we can write

$$\Delta T = \frac{\partial T}{\partial x} \Delta x + \frac{\partial T}{\partial y} \Delta y + \frac{\partial T}{\partial z} \Delta z. \quad (2.13)$$

The left side of Eq. (2.13) is a scalar. The right side is the sum of three products with Δx , Δy , and Δz , which are the components of a vector. It follows that the three numbers

$$\frac{\partial T}{\partial x}, \frac{\partial T}{\partial y}, \frac{\partial T}{\partial z}$$

are also the x -, y -, and z -components of a vector. We write this new vector with the symbol ∇T . The symbol ∇ (called “del”) is an upside-down Δ , and is supposed to remind us of differentiation. People read ∇T in various ways: “del- T ,” or “gradient of T ,” or “grad T ,”

$$\text{grad } T = \nabla T = \left(\frac{\partial T}{\partial x}, \frac{\partial T}{\partial y}, \frac{\partial T}{\partial z} \right)^*. \quad (2.14)$$

Using this notation, we can rewrite Eq. (2.13) in the more compact form

$$\Delta T = \nabla T \cdot \Delta \mathbf{R}. \quad (2.15)$$

In words, this equation says that the difference in temperature between two nearby points is the dot product of the gradient of T and the vector displacement between the points. The form of Eq. (2.15) also illustrates clearly our proof above that ∇T is indeed a vector.

Perhaps you are still not convinced? Let’s prove it in a different way. (Although if you look carefully, you may be able to see that it’s really the same proof in a longer-winded form!) We shall show that the components of ∇T transform in just the same way that components of \mathbf{R} do. If they do, ∇T is a vector according to our original definition of a vector in Chapter 11 of Vol. I. We take a new coordinate system x' , y' , z' , and in this new system we calculate $\partial T / \partial x'$, $\partial T / \partial y'$, and $\partial T / \partial z'$. To make things a little simpler, we let $z = z'$, so that we can forget about the z -coordinate. (You can check out the more general case for yourself.)

We take an $x'y'$ -system rotated an angle θ with respect to the xy -system, as in Fig. 2-6(a). For a point (x, y) the coordinates in the prime system are

$$x' = x \cos \theta + y \sin \theta, \quad (2.16)$$

$$y' = -x \sin \theta + y \cos \theta. \quad (2.17)$$

Or, solving for x and y ,

$$x = x' \cos \theta - y' \sin \theta, \quad (2.18)$$

$$y = x' \sin \theta + y' \cos \theta. \quad (2.19)$$

If any pair of numbers transforms with these equations in the same way that x and y do, they are the components of a vector.

Now let’s look at the difference in temperature between the two nearby points P_1 and P_2 , chosen as in Fig. 2-6(b). If we calculate with the x - and y -coordinates, we would write

$$\Delta T = \frac{\partial T}{\partial x} \Delta x \quad (2.20)$$

—since Δy is zero.

* In our notation, the expression (a, b, c) represents a vector with components a , b , and c . If you like to use the unit vectors i, j , and k , you may write

$$\nabla T = i \frac{\partial T}{\partial x} + j \frac{\partial T}{\partial y} + k \frac{\partial T}{\partial z}.$$

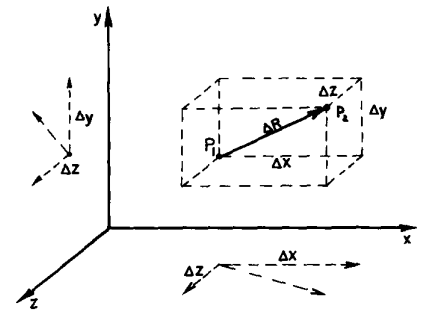


Fig. 2-5. The vector $\Delta \mathbf{R}$, whose components are Δx , Δy , and Δz .

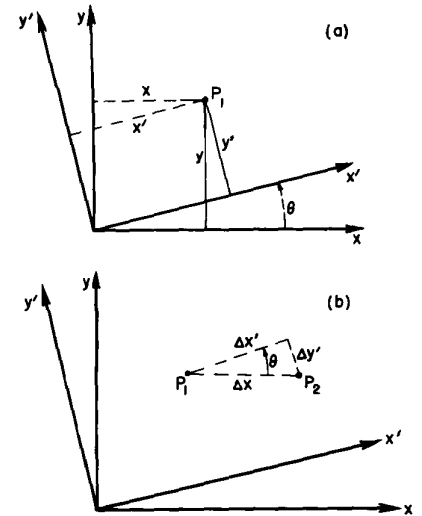


Fig. 2-6. (a) Transformation to a rotated coordinate system. (b) Special case of an interval $\Delta \mathbf{R}$ parallel to the x -axis.

What would a computation in the prime system give? We would have written

$$\Delta T = \frac{\partial T}{\partial x'} \Delta x' + \frac{\partial T}{\partial y'} \Delta y'. \quad (2.21)$$

Looking at Fig. 2-6(b), we see that

$$\Delta x' = \Delta x \cos \theta \quad (2.22)$$

and

$$\Delta y' = -\Delta x \sin \theta, \quad (2.23)$$

since Δy is negative when Δx is positive. Substituting these in Eq. (2.21), we find that

$$\Delta T = \frac{\partial T}{\partial x'} \Delta x \cos \theta - \frac{\partial T}{\partial y'} \Delta x \sin \theta \quad (2.24)$$

$$= \left(\frac{\partial T}{\partial x'} \cos \theta - \frac{\partial T}{\partial y'} \sin \theta \right) \Delta x. \quad (2.25)$$

Comparing Eq. (2.25) with (2.20), we see that

$$\frac{\partial T}{\partial x} = \frac{\partial T}{\partial x'} \cos \theta - \frac{\partial T}{\partial y'} \sin \theta. \quad (2.26)$$

This equation says that $\partial T/\partial x$ is obtained from $\partial T/\partial x'$ and $\partial T/\partial y'$, just as x is obtained from x' and y' in Eq. (2.18). So $\partial T/\partial x$ is the x -component of a vector. The same kind of arguments would show that $\partial T/\partial y$ and $\partial T/\partial z$ are y - and z -components. So ∇T is definitely a vector. It is a vector field derived from the scalar field T .

2-4 The operator ∇

Now we can do something that is extremely amusing and ingenious—and characteristic of the things that make mathematics beautiful. The argument that $\text{grad } T$, or ∇T , is a vector did not depend upon *what* scalar field we were differentiating. All the arguments would go the same if T were replaced by *any scalar field*. Since the transformation equations are the same no matter what we differentiate, we could just as well omit the T and replace Eq. (2.26) by the operator equation

$$\frac{\partial}{\partial x} = \frac{\partial}{\partial x'} \cos \theta - \frac{\partial}{\partial y'} \sin \theta. \quad (2.27)$$

We leave the operators, as Jeans said, “hungry for something to differentiate.”

Since the differential operators themselves transform as the components of a vector should, we can call them components of a *vector operator*. We can write

$$\nabla = \left(\frac{\partial}{\partial x}, \frac{\partial}{\partial y}, \frac{\partial}{\partial z} \right), \quad (2.28)$$

which means, of course,

$$\nabla_x = \frac{\partial}{\partial x}, \quad \nabla_y = \frac{\partial}{\partial y}, \quad \nabla_z = \frac{\partial}{\partial z}. \quad (2.29)$$

We have abstracted the gradient away from the T —that is the wonderful idea.

You must always remember, of course, that ∇ is an operator. Alone, it means nothing. If ∇ by itself means nothing, what does it mean if we multiply it by a scalar—say T —to get the product $T\nabla$? (One can always multiply a vector by a scalar.) It still does not mean anything. Its x -component is

$$T \frac{\partial}{\partial x}, \quad (2.30)$$

which is not a number, but is still some kind of operator. However, according to the algebra of vectors we would still call $T\nabla$ a vector.

Now let's multiply ∇ by a scalar on the other side, so that we have the product (∇T) . In ordinary algebra

$$TA = AT, \quad (2.31)$$

but we have to remember that operator algebra is a little different from ordinary vector algebra. With operators we must always keep the sequence right, so that the operations make proper sense. You will have no difficulty if you just remember that the operator ∇ obeys the same convention as the derivative notation. What is to be differentiated must be placed on the right of the ∇ . The order is important.

Keeping in mind this problem of order, we understand that $T\nabla$ is an operator, but the product ∇T is no longer a hungry operator; the operator is completely satisfied. It is indeed a physical vector having a meaning. It represents the spatial rate of change of T . The x -component of ∇T is how fast T changes in the x -direction. What is the direction of the vector ∇T ? We know that the rate of change of T in any direction is the component of ∇T in that direction (see Eq. 2.15). It follows that the direction of ∇T is that in which it has the largest possible component—in other words, the direction in which T changes the fastest. The gradient of T has the direction of the steepest uphill slope (in T).

2-5 Operations with ∇

Can we do any other algebra with the vector operator ∇ ? Let us try combining it with a vector. We can combine two vectors by making a dot product. We could make the products

$$(\text{a vector}) \cdot \nabla, \quad \text{or} \quad \nabla \cdot (\text{a vector}).$$

The first one doesn't mean anything yet, because it is still an operator. What it might ultimately mean would depend on what it is made to operate on. The second product is some scalar field. ($A \cdot B$ is always a scalar.)

Let's try the dot product of ∇ with a vector field we know, say \mathbf{h} . We write out the components:

$$\nabla \cdot \mathbf{h} = \nabla_x h_x + \nabla_y h_y + \nabla_z h_z \quad (2.32)$$

or

$$\nabla \cdot \mathbf{h} = \frac{\partial h_x}{\partial x} + \frac{\partial h_y}{\partial y} + \frac{\partial h_z}{\partial z}. \quad (2.33)$$

The sum is invariant under a coordinate transformation. If we were to choose a different system (indicated by primes), we would have*

$$\nabla' \cdot \mathbf{h} = \frac{\partial h_{x'}}{\partial x'} + \frac{\partial h_{y'}}{\partial y'} + \frac{\partial h_{z'}}{\partial z'}, \quad (2.34)$$

which is the *same* number as would be gotten from Eq. (2.33), even though it looks different. That is,

$$\nabla' \cdot \mathbf{h} = \nabla \cdot \mathbf{h} \quad (2.35)$$

for every point in space. So $\nabla \cdot \mathbf{h}$ is a scalar field, which must represent some physical quantity. You should realize that the combination of derivatives in $\nabla \cdot \mathbf{h}$ is rather special. There are all sorts of other combinations like $\partial h_y / \partial x$, which are neither scalars nor components of vectors.

The scalar quantity $\nabla \cdot (\text{a vector})$ is extremely useful in physics. It has been given the name the *divergence*. For example,

$$\nabla \cdot \mathbf{h} = \text{div } \mathbf{h} = \text{"divergence of } \mathbf{h} \text{"} \quad (2.36)$$

As we did for ∇T , we can ascribe a physical significance to $\nabla \cdot \mathbf{h}$. We shall, however, postpone that until later.

* We think of \mathbf{h} as a *physical* quantity that depends on position in space, and not strictly as a mathematical function of three variables. When \mathbf{h} is "differentiated" with respect to x , y , and z , or with respect to x' , y' , and z' , the mathematical expression for \mathbf{h} must first be expressed as a function of the appropriate variables.

First, we wish to see what else we can cook up with the vector operator ∇ . What about a cross product? We must expect that

$$\nabla \times \mathbf{h} = \text{a vector.} \quad (2.37)$$

It is a vector whose components we can write by the usual rule for cross products (see Eq. 2.2):

$$(\nabla \times \mathbf{h})_z = \nabla_x h_y - \nabla_y h_x = \frac{\partial h_y}{\partial x} - \frac{\partial h_x}{\partial y}. \quad (2.38)$$

Similarly,

$$(\nabla \times \mathbf{h})_x = \nabla_y h_z - \nabla_z h_y = \frac{\partial h_z}{\partial y} - \frac{\partial h_y}{\partial z} \quad (2.39)$$

and

$$(\nabla \times \mathbf{h})_y = \nabla_z h_x - \nabla_x h_z = \frac{\partial h_x}{\partial z} - \frac{\partial h_z}{\partial x}. \quad (2.40)$$

The combination $\nabla \times \mathbf{h}$ is called “the *curl* of \mathbf{h} .” The reason for the name and the physical meaning of the combination will be discussed later.

Summarizing, we have three kinds of combinations with ∇ :

$$\nabla T = \text{grad } T = \text{a vector,}$$

$$\nabla \cdot \mathbf{h} = \text{div } \mathbf{h} = \text{a scalar,}$$

$$\nabla \times \mathbf{h} = \text{curl } \mathbf{h} = \text{a vector.}$$

Using these combinations, we can write about the spatial variations of fields in a convenient way—in a way that is general, in that it doesn’t depend on any particular set of axes.

As an example of the use of our vector differential operator ∇ , we write a set of vector equations which contain the same laws of electromagnetism that we gave in words in Chapter 1. They are called Maxwell’s equations.

Maxwell’s Equations

$$\begin{aligned} (1) \quad \nabla \cdot \mathbf{E} &= \frac{\rho}{\epsilon_0} \\ (2) \quad \nabla \times \mathbf{E} &= -\frac{\partial \mathbf{B}}{\partial t} \\ (3) \quad \nabla \cdot \mathbf{B} &= 0 \\ (4) \quad c^2 \nabla \times \mathbf{B} &= \frac{\partial \mathbf{E}}{\partial t} + \frac{\mathbf{j}}{\epsilon_0} \end{aligned} \quad (2.41)$$

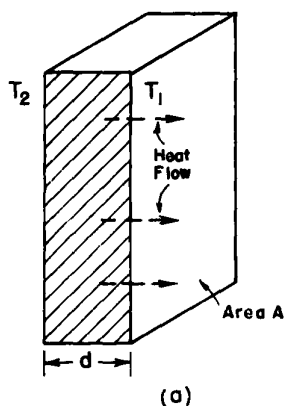
where ρ (rho), the “electric charge density,” is the amount of charge per unit volume, and \mathbf{j} , the “electric current density,” is the rate at which charge flows through a unit area per second. These four equations contain the complete classical theory of the electromagnetic field. You see what an elegantly simple form we can get with our new notation!

2-6 The differential equation of heat flow

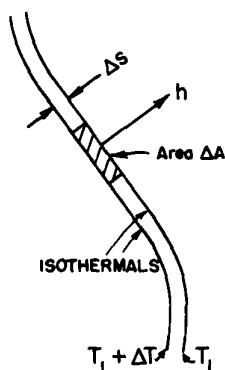
Let us give another example of a law of physics written in vector notation. The law is not a precise one, but for many metals and a number of other substances that conduct heat it is quite accurate. You know that if you take a slab of material and heat one face to temperature T_2 and cool the other to a different temperature T_1 , the heat will flow through the material from T_2 to T_1 [Fig. 2-7(a)]. The heat flow is proportional to the area A of the faces, and to the temperature difference. It is also inversely proportional to d , the distance between the plates. (For a given temperature difference, the thinner the slab the greater the heat flow.) Letting J be the thermal energy that passes per unit time through the slab, we write

$$J = \kappa(T_2 - T_1) \frac{A}{d}. \quad (2.42)$$

The constant of proportionality κ (kappa) is called the *thermal conductivity*.



(a)



(b)

Fig. 2-7. (a) Heat flow through a slab. (b) An infinitesimal slab parallel to an isothermal surface in a large block.

What will happen in a more complicated case? Say in an odd-shaped block of material in which the temperature varies in peculiar ways? Suppose we look at a tiny piece of the block and imagine a slab like that of Fig. 2-7(a) on a miniature scale. We orient the faces parallel to the isothermal surfaces, as in Fig. 2-7(b), so that Eq. (2.42) is correct for the small slab.

If the area of the small slab is ΔA , the heat flow per unit time is

$$\Delta J = \kappa \Delta T \frac{\Delta A}{\Delta s}, \quad (2.43)$$

where Δs is the thickness of the slab. Now $\Delta J/\Delta A$ we have defined earlier as the magnitude of \mathbf{h} , whose direction is the heat flow. The heat flow will be from $T_1 + \Delta T$ toward T_1 , and so it will be perpendicular to the isotherms, as drawn in Fig. 2-7(b). Also, $\Delta T/\Delta s$ is just the rate of change of T with position. And since the position change is perpendicular to the isotherms, our $\Delta T/\Delta s$ is the maximum rate of change. It is, therefore, just the magnitude of ∇T . Now since the direction of ∇T is opposite to that of \mathbf{h} , we can write (2.43) as a vector equation:

$$\mathbf{h} = -\kappa \nabla T. \quad (2.44)$$

(The minus sign is necessary because heat flows “downhill” in temperature.) Equation (2.44) is the differential equation of heat conduction in bulk materials. You see that it is a proper vector equation. Each side is a vector if κ is just a number. It is the generalization to arbitrary cases of the special relation (2.42) for rectangular slabs. Later we should learn to write all sorts of elementary physics relations like (2.42) in the more sophisticated vector notation. This notation is useful not only because it makes the equations *look* simpler. It also shows most clearly the *physical content* of the equations without reference to any arbitrarily chosen coordinate system.

2-7 Second derivatives of vector fields

So far we have had only first derivatives. Why not second derivatives? We could have several combinations:

$$\begin{aligned} (a) \quad & \nabla \cdot (\nabla T) \\ (b) \quad & \nabla \times (\nabla T) \\ (c) \quad & \nabla(\nabla \cdot \mathbf{h}) \\ (d) \quad & \nabla \cdot (\nabla \times \mathbf{h}) \\ (e) \quad & \nabla \times (\nabla \times \mathbf{h}) \end{aligned} \quad (2.45)$$

You can check that these are all the possible combinations.

Let's look first at the second one, (b). It has the same form as

$$\mathbf{A} \times (\mathbf{A}T) = (\mathbf{A} \times \mathbf{A})T = 0,$$

since $\mathbf{A} \times \mathbf{A}$ is always zero. So we should have

$$\text{curl}(\text{grad } T) = \nabla \times (\nabla T) = 0. \quad (2.46)$$

We can see how this equation comes about if we go through once with the components:

$$\begin{aligned} [\nabla \times (\nabla T)]_z &= \nabla_x(\nabla T)_y - \nabla_y(\nabla T)_x \\ &= \frac{\partial}{\partial x} \left(\frac{\partial T}{\partial y} \right) - \frac{\partial}{\partial y} \left(\frac{\partial T}{\partial x} \right), \end{aligned} \quad (2.47)$$

which is zero (by Eq. 2.8). It goes the same for the other components. So $\nabla \times (\nabla T) = 0$, for any temperature distribution—in fact, for *any* scalar function.

Now let us take another example. Let us see whether we can find another zero. The dot product of a vector with a cross product which contains that vector is zero:

$$\mathbf{A} \cdot (\mathbf{A} \times \mathbf{B}) = 0. \quad (2.48)$$

because $\mathbf{A} \times \mathbf{B}$ is perpendicular to \mathbf{A} , and so has no components in the direction \mathbf{A} . The same combination appears in (d) of (2.45), so we have

$$\nabla \cdot (\nabla \times \mathbf{h}) = \text{div}(\text{curl } \mathbf{h}) = 0. \quad (2.49)$$

Again, it is easy to show that it is zero by carrying through the operations with components.

Now we are going to state two mathematical theorems that we will not prove. They are very interesting and useful theorems for physicists to know.

In a physical problem we frequently find that the curl of some quantity—say of the vector field \mathbf{A} —is zero. Now we have seen (Eq. 2.46) that the curl of a gradient is zero, which is easy to remember because of the way the vectors work. It could certainly be, then, that \mathbf{A} is the gradient of some quantity, because then its curl would necessarily be zero. The interesting theorem is that if the curl \mathbf{A} is zero, then \mathbf{A} is *always* the gradient of *something*—there is some scalar field ψ (psi) such that \mathbf{A} is equal to $\text{grad } \psi$. In other words, we have the

THEOREM:

$$\begin{array}{ll} \text{If} & \nabla \times \mathbf{A} = 0 \\ \text{there is a} & \psi \\ \text{such that} & \mathbf{A} = \nabla \psi. \end{array} \quad (2.50)$$

There is a similar theorem if the divergence of \mathbf{A} is zero. We have seen in Eq. (2.49) that the divergence of a curl of something is always zero. If you come across a vector field \mathbf{D} for which $\text{div } \mathbf{D}$ is zero, then you can conclude that \mathbf{D} is the curl of some vector field \mathbf{C} .

THEOREM:

$$\begin{array}{ll} \text{If} & \nabla \cdot \mathbf{D} = 0 \\ \text{there is a} & \mathbf{C} \\ \text{such that} & \mathbf{D} = \nabla \times \mathbf{C}. \end{array} \quad (2.51)$$

In looking at the possible combinations of two ∇ operators, we have found that two of them always give zero. Now we look at the ones that are *not* zero. Take the combination $\nabla \cdot (\nabla T)$, which was first on our list. It is not, in general, zero. We write out the components:

$$\nabla T = \nabla_x T + \nabla_y T + \nabla_z T.$$

Then

$$\begin{aligned} \nabla \cdot (\nabla T) &= \nabla_x(\nabla_x T) + \nabla_y(\nabla_y T) + \nabla_z(\nabla_z T) \\ &= \frac{\partial^2 T}{\partial x^2} + \frac{\partial^2 T}{\partial y^2} + \frac{\partial^2 T}{\partial z^2}, \end{aligned} \quad (2.52)$$

which would, in general, come out to be some number. It is a scalar field.

You see that we do not need to keep the parentheses, but can write, without any chance of confusion,

$$\nabla \cdot (\nabla T) = \nabla \cdot \nabla T = (\nabla \cdot \nabla)T = \nabla^2 T. \quad (2.53)$$

We look at ∇^2 as a new operator. It is a scalar operator. Because it appears often in physics, it has been given a special name—the *Laplacian*.

$$\text{Laplacian} = \nabla^2 = \frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} + \frac{\partial^2}{\partial z^2}. \quad (2.54)$$

Since the Laplacian is a scalar operator, we may operate with it on a vector—by which we mean the same operation on each component in rectangular coordinates:

$$\nabla^2 \mathbf{h} = (\nabla^2 h_x, \nabla^2 h_y, \nabla^2 h_z).$$

Let's look at one more possibility: $\nabla \times (\nabla \times \mathbf{h})$, which was (e) in the list (2.45). Now the curl of the curl can be written differently if we use the vector equality (2.6):

$$\mathbf{A} \times (\mathbf{B} \times \mathbf{C}) = \mathbf{B}(\mathbf{A} \cdot \mathbf{C}) - \mathbf{C}(\mathbf{A} \cdot \mathbf{B}). \quad (2.55)$$

In order to use this formula, we should replace \mathbf{A} and \mathbf{B} by the operator ∇ and put $\mathbf{C} = \mathbf{h}$. If we do that, we get

$$\nabla \times (\nabla \times \mathbf{h}) = \nabla(\nabla \cdot \mathbf{h}) - \mathbf{h}(\nabla \cdot \nabla) \dots ???$$

Wait a minute! Something is wrong. The first two terms are vectors all right (the operators are satisfied), but the last term doesn't come out to anything. It's still an operator. The trouble is that we haven't been careful enough about keeping the order of our terms straight. If you look again at Eq. (2.55), however, you see that we could equally well have written it as

$$\mathbf{A} \times (\mathbf{B} \times \mathbf{C}) = \mathbf{B}(\mathbf{A} \cdot \mathbf{C}) - (\mathbf{A} \cdot \mathbf{B})\mathbf{C}. \quad (2.56)$$

The order of terms looks better. Now let's make our substitution in (2.56). We get

$$\nabla \times (\nabla \times \mathbf{h}) = \nabla(\nabla \cdot \mathbf{h}) - (\nabla \cdot \nabla)\mathbf{h}. \quad (2.57)$$

This form looks all right. It is, in fact, correct, as you can verify by computing the components. The last term is the Laplacian, so we can equally well write

$$\nabla \times (\nabla \times \mathbf{h}) = \nabla(\nabla \cdot \mathbf{h}) - \nabla^2 \mathbf{h}. \quad (2.58)$$

We have had something to say about all of the combinations in our list of double ∇ 's, except for (c), $\nabla(\nabla \cdot \mathbf{h})$. It is a possible vector field, but there is nothing special to say about it. It's just some vector field which may occasionally come up.

It will be convenient to have a table of our conclusions:

- (a) $\nabla \cdot (\nabla T) = \nabla^2 T =$ a scalar field
- (b) $\nabla \times (\nabla T) = 0$
- (c) $\nabla(\nabla \cdot \mathbf{h}) =$ a vector field
- (d) $\nabla \cdot (\nabla \times \mathbf{h}) = 0$
- (e) $\nabla \times (\nabla \times \mathbf{h}) = \nabla(\nabla \cdot \mathbf{h}) - \nabla^2 \mathbf{h}$
- (f) $(\nabla \cdot \nabla)\mathbf{h} = \nabla^2 \mathbf{h} =$ a vector field

You may notice that we haven't tried to invent a new vector operator ($\nabla \times \nabla$). Do you see why?

2-8 Pitfalls

We have been applying our knowledge of ordinary vector algebra to the algebra of the operator ∇ . We have to be careful, though, because it is possible to go astray. There are two pitfalls which we will mention, although they will not come up in this course. What would you say about the following expression, that involves the two scalar functions ψ and ϕ (phi):

$$(\nabla\psi) \times (\nabla\phi)?$$

You might want to say: it must be zero because it's just like

$$(\mathbf{A}\mathbf{a}) \times (\mathbf{A}\mathbf{b}),$$

which is zero because the cross product of two *equal* vectors $A \times A$ is always zero. But in our example the two operators ∇ are not equal! The first one operates on one function, ψ ; the other operates on a different function, ϕ . So although we represent them by the same symbol ∇ , they must be considered as different operators. Clearly, the direction of $\nabla\psi$ depends on the function ψ , so it is not likely to be parallel to $\nabla\phi$.

$$(\nabla\psi) \times (\nabla\phi) \neq 0 \quad (\text{generally}).$$

Fortunately, we won't have to use such expressions. (What we have said doesn't change the fact that $\nabla \times \nabla\psi = 0$ for any scalar field, because here both ∇ 's operate on the same function.)

Pitfall number two (which, again, we need not get into in our course) is the following: The rules that we have outlined here are simple and nice when we use rectangular coordinates. For example, if we have $\nabla^2 h$ and we want the x -component, it is

$$(\nabla^2 h)_x = \left(\frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} + \frac{\partial^2}{\partial z^2} \right) h_x = \nabla^2 h_x. \quad (2.60)$$

The same expression would *not* work if we were to ask for the *radial* component of $\nabla^2 h$. The radial component of $\nabla^2 h$ is not equal to $\nabla^2 h_r$. The reason is that when we are dealing with the algebra of vectors, the directions of the vectors are all quite definite. But when we are dealing with vector fields, their directions are different at different places. If we try to describe a vector field in, say, polar coordinates, what we call the "radial" direction varies from point to point. So we can get into a lot of trouble when we start to differentiate the components. For example, even for a *constant* vector field, the radial component changes from point to point.

It is usually safest and simplest just to stick to rectangular coordinates and avoid trouble, but there is one exception worth mentioning: Since the Laplacian ∇^2 , is a scalar, we can write it in any coordinate system we want to (for example, in polar coordinates). But since it is a differential operator, we should use it only on vectors whose components are in a fixed direction—that means rectangular coordinates. So we shall express all of our vector fields in terms of their x -, y -, and z -components when we write our vector differential equations out in components.

Vector Integral Calculus

3-1 Vector integrals; the line integral of $\nabla\psi$

We found in Chapter 2 that there were various ways of taking derivatives of fields. Some gave vector fields; some gave scalar fields. Although we developed many different formulas, everything in Chapter 2 could be summarized in one rule: the operators $\partial/\partial x$, $\partial/\partial y$, and $\partial/\partial z$ are the three components of a vector operator ∇ . We would now like to get some understanding of the significance of the derivatives of fields. We will then have a better feeling for what a vector field equation means.

We have already discussed the meaning of the gradient operation (∇ on a scalar). Now we turn to the meanings of the divergence and curl operations. The interpretation of these quantities is best done in terms of certain vector integrals and equations relating such integrals. These equations cannot, unfortunately, be obtained from vector algebra by some easy substitution, so you will just have to learn them as something new. Of these integral formulas, one is practically trivial, but the other two are not. We will derive them and explain their implications. The equations we shall study are really mathematical theorems. They will be useful not only for interpreting the meaning and the content of the divergence and the curl, but also in working out general physical theories. These mathematical theorems are, for the theory of fields, what the theorem of the conservation of energy is to the mechanics of particles. General theorems like these are important for a deeper understanding of physics. You will find, though, that they are not very useful for solving problems—except in the simplest cases. It is delightful, however, that in the beginning of our subject there will be many simple problems which can be solved with the three integral formulas we are going to treat. We will see, however, as the problems get harder, that we can no longer use these simple methods.

We take up first an integral formula involving the gradient. The relation contains a very simple idea: Since the gradient represents the rate of change of a field quantity, if we integrate that rate of change, we should get the total change. Suppose we have the scalar field $\psi(x, y, z)$. At any two points (1) and (2), the function ψ will have the values $\psi(1)$ and $\psi(2)$, respectively. [We use a convenient notation, in which (2) represents the point (x_2, y_2, z_2) and $\psi(2)$ means the same thing as $\psi(x_2, y_2, z_2)$.] If Γ (gamma) is any curve joining (1) and (2), as in Fig. 3-1, the following relation is true:

THEOREM 1.

$$\psi(2) - \psi(1) = \int_{\text{along } \Gamma}^{(2)} (\nabla\psi) \cdot ds. \quad (3.1)$$

The integral is a *line integral*, from (1) to (2) along the curve Γ , of the dot product of $\nabla\psi$ —a vector—with ds —another vector which is an infinitesimal line element of the curve Γ (directed away from (1) and toward (2)).

First, we should review what we mean by a line integral. Consider a scalar function $f(x, y, z)$, and the curve Γ joining two points (1) and (2). We mark off the curve at a number of points and join these points by straight-line segments, as shown in Fig. 3-2. Each segment has the length Δs_i , where i is an index that runs 1, 2, 3, . . . By the line integral

$$\int_{\text{along } \Gamma}^{(2)} f ds$$

3-1 Vector integrals; the line integral of $\nabla\psi$

3-2 The flux of a vector field

3-3 The flux from a cube; Gauss' theorem

3-4 Heat conduction; the diffusion equation

3-5 The circulation of a vector field

3-6 The circulation around a square; Stokes' theorem

3-7 Curl-free and divergence-free fields

3-8 Summary

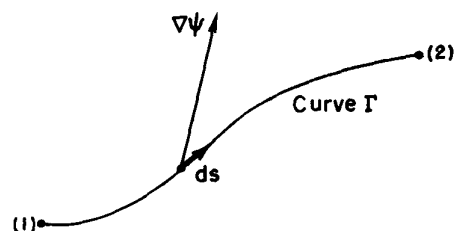


Fig. 3-1. The terms used in Eq. (3.1). The vector $\nabla\psi$ is evaluated at the line element ds .

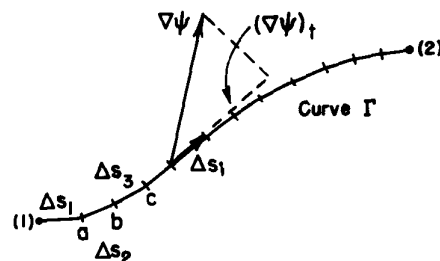


Fig. 3-2. The line integral is the limit of a sum.

we mean the limit of the sum

$$\sum f_i \Delta s_i,$$

where f_i is the value of the function at the i th segment. The limiting value is what the sum approaches as we add more and more segments (in a sensible way, so that the largest $\Delta s_i \rightarrow 0$).

The integral in our theorem, Eq. (3.1), means the same thing, although it looks a little different. Instead of f , we have another scalar—the component of $\nabla\psi$ in the direction of Δs . If we write $(\nabla\psi)_t$ for this tangential component, it is clear that

$$(\nabla\psi)_t \Delta s = (\nabla\psi) \cdot \Delta s. \quad (3.2)$$

The integral in Eq. (3.1) means the sum of such terms.

Now let's see why Eq. (3.1) is true. In Chapter 1, we showed that the component of $\nabla\psi$ along a small displacement ΔR was the rate of change of ψ in the direction of ΔR . Consider the line segment Δs from (1) to point a in Fig. 3-2. According to our definition,

$$\Delta\psi_1 = \psi(a) - \psi(1) = (\nabla\psi)_1 \cdot \Delta s_1. \quad (3.3)$$

Also, we have

$$\psi(b) - \psi(a) = (\nabla\psi)_2 \cdot \Delta s_2, \quad (3.4)$$

where, of course, $(\nabla\psi)_1$ means the gradient evaluated at the segment Δs_1 , and $(\nabla\psi)_2$, the gradient evaluated at Δs_2 . If we add Eqs. (3.3) and (3.4), we get

$$\psi(b) - \psi(1) = (\nabla\psi)_1 \cdot \Delta s_1 + (\nabla\psi)_2 \cdot \Delta s_2. \quad (3.5)$$

You can see that if we keep adding such terms, we get the result

$$\psi(2) - \psi(1) = \sum (\nabla\psi)_i \cdot \Delta s_i. \quad (3.6)$$

The left-hand side doesn't depend on how we choose our intervals—if (1) and (2) are kept always the same—so we can take the limit of the right-hand side. We have therefore proved Eq. (3.1).

You can see from our proof that just as the equality doesn't depend on how the points a, b, c, \dots , are chosen, similarly it doesn't depend on what we choose for the curve Γ to join (1) and (2). Our theorem is correct for *any* curve from (1) to (2).

One remark on notation: You will see that there is no confusion if we write, for convenience,

$$(\nabla\psi) \cdot ds = \nabla\psi \cdot ds. \quad (3.7)$$

With this notation, our theorem is

$$\text{THEOREM 1.} \quad \psi(2) - \psi(1) = \int_{\text{any curve from (1) to (2)}}^{(2)} \nabla\psi \cdot ds. \quad (3.8)$$

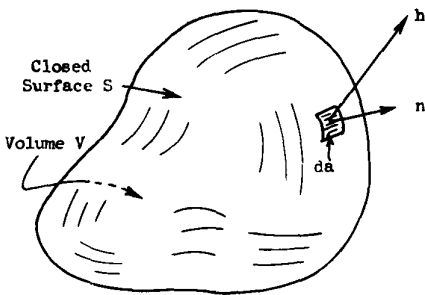


Fig. 3-3. The closed surface S defines the volume V . The unit vector n is the outward normal to the surface element da , and h is the heat-flow vector at the surface element.

3-2 The flux of a vector field

Before we consider our next integral theorem—a theorem about the divergence—we would like to study a certain idea which has an easily understood physical significance in the case of heat flow. We have defined the vector h , which represents the heat that flows through a unit area in a unit time. Suppose that inside a block of material we have some closed surface S which encloses the volume V (Fig. 3-3). We would like to find out how much heat is flowing out of this volume. We can, of course, find it by calculating the total heat flow out of the surface S .

We write da for the area of an element of the surface. The symbol stands for a two-dimensional differential. If, for instance, the area happened to be in the xy -plane we would have

$$da = dx dy.$$

Later we shall have integrals over volume and for these it is convenient to consider a differential volume that is a little cube. So when we write dV we mean

$$dV = dx dy dz.$$

Some people like to write d^2a instead of da to remind themselves that it is kind of a second-order quantity. They would also write d^3V instead of dV . We will use the simpler notation, and assume that you can remember that an area has two dimensions and a volume has three.

The heat flow out through the surface element da is the area times the component of \mathbf{h} perpendicular to da . We have already defined \mathbf{n} as a unit vector pointing outward at right angles to the surface (Fig. 3-3). The component of \mathbf{h} that we want is

$$h_n = \mathbf{h} \cdot \mathbf{n}. \quad (3.9)$$

The heat flow out through da is then

$$\mathbf{h} \cdot \mathbf{n} da. \quad (3.10)$$

To get the total heat flow through any surface we sum the contributions from all the elements of the surface. In other words, we integrate (3.10) over the whole surface:

$$\text{Total heat flow outward through } S = \int_S \mathbf{h} \cdot \mathbf{n} da. \quad (3.11)$$

We are also going to call this surface integral “the flux of \mathbf{h} through the surface.” Originally the word flux meant flow, so that the surface integral just means the flow of \mathbf{h} through the surface. We may think: \mathbf{h} is the “current density” of heat flow and the surface integral of it is the total heat current directed out of the surface; that is, the thermal energy per unit time (joules per second).

We would like to generalize this idea to the case where the vector does not represent the flow of anything; for instance, it might be the electric field. We can certainly still integrate the normal component of the electric field over an area if we wish. Although it is not the flow of anything, we still call it the “flux.” We say

$$\text{Flux of } \mathbf{E} \text{ through the surface } S = \int_S \mathbf{E} \cdot \mathbf{n} da. \quad (3.12)$$

We generalize the word “flux” to mean the “surface integral of the normal component” of a vector. We will also use the same definition even when the surface considered is not a closed one, as it is here.

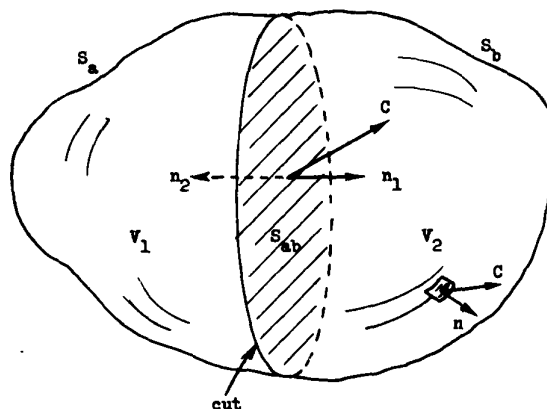
Returning to the special case of heat flow, let us take a situation in which *heat is conserved*. For example, imagine some material in which after an initial heating no further heat energy is generated or absorbed. Then, if there is a net heat flow out of a closed surface, the heat content of the volume inside must decrease. So, in circumstances in which heat would be conserved, we say that

$$\int_S \mathbf{h} \cdot \mathbf{n} da = - \frac{dQ}{dt}, \quad (3.13)$$

where Q is the heat inside the surface. The heat flux out of S is equal to minus the rate of change with respect to time of the total heat Q inside of S . This interpretation is possible because we are speaking of heat flow and also because we supposed that the heat was conserved. We could not, of course, speak of the total heat inside the volume if heat were being generated there.

Now we shall point out an interesting fact about the flux of any vector. You may think of the heat flow vector if you wish, but what we say will be true for any vector field \mathbf{C} . Imagine that we have a closed surface S that encloses the volume V . We now separate the volume into two parts by some kind of a “cut,” as in Fig. 3-4. Now we have two closed surfaces and volumes. The volume V_1 is enclosed in the surface S_1 , which is made up of part of the original surface S_a and of the surface of the cut, S_{ab} . The volume V_2 is enclosed by S_2 , which is made up of the rest of the original surface S_b and closed off by the cut S_{ab} . Now consider the

Fig. 3-4. A volume V contained inside the surface S is divided into two pieces by a "cut" at the surface S_{ab} . We now have the volume V_1 enclosed in the surface $S_1 = S_a + S_{ab}$ and the volume V_2 enclosed in the surface $S_2 = S_b + S_{ab}$.



following question: Suppose we calculate the flux out through surface S_1 and add to it the flux through surface S_2 . Does the sum equal the flux through the whole surface that we started with? The answer is yes. The flux through the part of the surfaces S_{ab} common to both S_1 and S_2 just exactly cancels out. For the flux of the vector C out of V_1 , we can write

$$\text{Flux through } S_1 = \int_{S_a} C \cdot n \, da + \int_{S_{ab}} C \cdot n_1 \, da, \quad (3.14)$$

and for the flux out of V_2 ,

$$\text{Flux through } S_2 = \int_{S_b} C \cdot n \, da + \int_{S_{ab}} C \cdot n_2 \, da. \quad (3.15)$$

Note that in the second integral we have written n_1 for the outward normal for S_{ab} when it belongs to S_1 , and n_2 when it belongs to S_2 , as shown in Fig. 3-4. Clearly, $n_1 = -n_2$, so that

$$\int_{S_{ab}} C \cdot n_1 \, da = - \int_{S_{ab}} C \cdot n_2 \, da. \quad (3.16)$$

If we now add Eqs. (3.14) and (3.15), we see that the sum of the fluxes through S_1 and S_2 is just the sum of two integrals which, taken together, give the flux through the original surface $S = S_a + S_b$.

We see that the flux through the complete outer surface S can be considered as the sum of the fluxes from the two pieces into which the volume was broken. We can similarly subdivide again—say by cutting V_1 into two pieces. You see that the same arguments apply. So for *any* way of dividing the original volume, it must be generally true that the flux through the outer surface, which is the original integral, is equal to a sum of the fluxes out of all the little interior pieces.

3-3 The flux from a cube; Gauss' theorem

We now take the special case of a small cube* and find an interesting formula for the flux out of it. Consider a cube whose edges are lined up with the axes as in Fig. 3-5. Let us suppose that the coordinates of the corner nearest the origin are x, y, z . Let Δx be the length of the cube in the x -direction, Δy be the length in the y -direction, and Δz be the length in the z -direction. We wish to find the flux of a vector field C through the surface of the cube. We shall do this by making a sum of the fluxes through each of the six faces. First, consider the face marked 1 in the figure. The flux *outward* on this face is the negative of the x -component of C , integrated over the area of the face. This flux is

$$- \int C_x \, dy \, dz.$$

Since we are considering a *small* cube, we can approximate this integral by the

* The following development applies equally well to any rectangular parallelepiped.

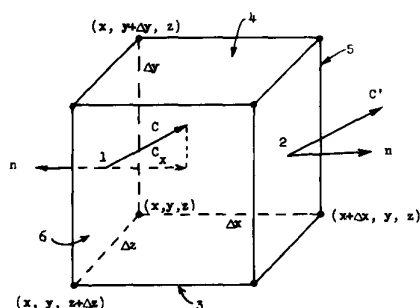


Fig. 3-5. Computation of the flux of C out of a small cube.

value of C_x at the center of the face—which we call the point (1)—multiplied by the area of the face, $\Delta y \Delta z$:

$$\text{Flux out of 1} = -C_x(1) \Delta y \Delta z.$$

Similarly, for the flux out of face 2, we write

$$\text{Flux out of 2} = C_x(2) \Delta y \Delta z.$$

Now $C_x(1)$ and $C_x(2)$ are, in general, slightly different. If Δx is small enough, we can write

$$C_x(2) = C_x(1) + \frac{\partial C_x}{\partial x} \Delta x.$$

There are, of course, more terms, but they will involve $(\Delta x)^2$ and higher powers, and so will be negligible if we consider only the limit of small Δx . So the flux through face 2 is

$$\text{Flux out of 2} = \left[C_x(1) + \frac{\partial C_x}{\partial x} \Delta x \right] \Delta y \Delta z.$$

Summing the fluxes for faces 1 and 2, we get

$$\text{Flux out of 1 and 2} = \frac{\partial C_x}{\partial x} \Delta x \Delta y \Delta z.$$

The derivative should really be evaluated at the center of face 1; that is, at $[x, y + (\Delta y/2), z + (\Delta z/2)]$. But in the limit of an infinitesimal cube, we make a negligible error if we evaluate it at the corner (x, y, z) .

Applying the same reasoning to each of the other pairs of faces, we have

$$\text{Flux out of 3 and 4} = \frac{\partial C_y}{\partial y} \Delta x \Delta y \Delta z$$

and

$$\text{Flux out of 5 and 6} = \frac{\partial C_z}{\partial z} \Delta x \Delta y \Delta z.$$

The total flux through all the faces is the sum of these terms. We find that

$$\int_{\text{cube}} \mathbf{C} \cdot \mathbf{n} \, da = \left(\frac{\partial C_x}{\partial x} + \frac{\partial C_y}{\partial y} + \frac{\partial C_z}{\partial z} \right) \Delta x \Delta y \Delta z,$$

and the sum of the derivatives is just $\nabla \cdot \mathbf{C}$. Also, $\Delta x \Delta y \Delta z = \Delta V$, the volume of the cube. So we can say that *for an infinitesimal cube*

$$\int_{\text{surface}} \mathbf{C} \cdot \mathbf{n} \, da = (\nabla \cdot \mathbf{C}) \Delta V. \quad (3.17)$$

We have shown that the outward flux from the surface of an infinitesimal cube is equal to the divergence of the vector multiplied by the volume of the cube. We now see the “meaning” of the divergence of a vector. The divergence of a vector at the point P is the flux—the outgoing “flow” of \mathbf{C} —*per unit volume*, in the neighborhood of P .

We have connected the divergence of \mathbf{C} to the flux of \mathbf{C} out of each infinitesimal volume. For any finite volume we can use the fact we proved above—that the total flux from a volume is the sum of the fluxes out of each part. We can, that is, integrate the divergence over the entire volume. This gives us the theorem that the integral of the normal component of any vector over any closed surface can also be written as the integral of the divergence of the vector over the volume enclosed by the surface. This theorem is named after Gauss.

GAUSS’ THEOREM.

$$\int_S \mathbf{C} \cdot \mathbf{n} \, da = \int_V \nabla \cdot \mathbf{C} \, dV, \quad (3.18)$$

where S is any closed surface and V is the volume inside it.

3-4 Heat conduction; the diffusion equation

Let's consider an example of the use of this theorem, just to get familiar with it. Suppose we take again the case of heat flow in, say, a metal. Suppose we have a simple situation in which all the heat has been previously put in and the body is just cooling off. There are no sources of heat, so that heat is conserved. Then how much heat is there inside some chosen volume at any time? It must be *decreasing* by just the amount that flows out of the surface of the volume. If our volume is a little cube, we would write, following Eq. (3.17),

$$\text{Heat out} = \int_{\text{cube}} \mathbf{h} \cdot \mathbf{n} \, da = \nabla \cdot \mathbf{h} \, \Delta V. \quad (3.19)$$

But this must equal the rate of loss of the heat inside the cube. If q is the heat per unit volume, the heat in the cube is $q \, \Delta V$, and the rate of *loss* is

$$- \frac{d}{dt} (q \, \Delta V) = - \frac{dq}{dt} \, \Delta V. \quad (3.20)$$

Comparing (3.19) and (3.20), we see that

$$- \frac{dq}{dt} = \nabla \cdot \mathbf{h}. \quad (3.21)$$

Take careful note of the form of this equation; the form appears often in physics. It expresses a conservation law—here the conservation of heat. We have expressed the same physical fact in another way in Eq. (3.13). Here we have the *differential* form of a conservation equation, while Eq. (3.13) is the *integral* form.

We have obtained Eq. (3.21) by applying Eq. (3.13) to an infinitesimal cube. We can also go the other way. For a big volume V bounded by S , Gauss' law says that

$$\int_S \mathbf{h} \cdot \mathbf{n} \, da = \int_V \nabla \cdot \mathbf{h} \, dV. \quad (3.22)$$

Using (3.21), the integral on the right-hand side is found to be just $-dQ/dt$, and again we have Eq. (3.13).

Now let's consider a different case. Imagine that we have a block of material and that inside it there is a very tiny hole in which some chemical reaction is taking place and generating heat. Or we could imagine that there are some wires running into a tiny resistor that is being heated by an electric current. We shall suppose that the heat is generated practically at a point, and let W represent the energy liberated per second at that point. We shall suppose that in the rest of the volume heat is conserved, and that the heat generation has been going on for a long time—so that now the temperature is no longer changing anywhere. The problem is: What does the heat vector \mathbf{h} look like at various places in the metal? How much heat flow is there at each point?

We know that if we integrate the normal component of \mathbf{h} over a closed surface that encloses the source, we will always get W . All the heat that is being generated at the point source must flow out through the surface, since we have supposed that the flow is steady. We have the difficult problem of finding a vector field which, when integrated over any surface, always gives W . We can, however, find the field rather easily by taking a somewhat special surface. We take a sphere of radius R , centered at the source, and assume that the heat flow is radial (Fig. 3-6). Our intuition tells us that \mathbf{h} should be radial if the block of material is large and we don't get too close to the edges, and it should also have the same magnitude at all points on the sphere. You see that we are adding a certain amount of guesswork—usually called “physical intuition”—to our mathematics in order to find the answer.

When \mathbf{h} is radial and spherically symmetric, the integral of the normal component of \mathbf{h} over the area is very simple, because the normal component is just

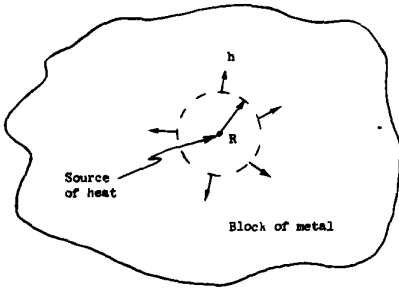


Fig. 3-6. In the region near a point source of heat, the heat flow is radially outward.

the magnitude of \mathbf{h} and is constant. The area over which we integrate is $4\pi R^2$. We have then that

$$\int_S \mathbf{h} \cdot \mathbf{n} \, da = h \cdot 4\pi R^2 \quad (3.23)$$

(where h is the magnitude of \mathbf{h}). This integral should equal W , the rate at which heat is produced at the source. We get

$$h = \frac{W}{4\pi R^2},$$

or

$$\mathbf{h} = \frac{W}{4\pi R^2} \mathbf{e}_r, \quad (3.24)$$

where, as usual, \mathbf{e}_r represents a unit vector in the radial direction. Our result says that \mathbf{h} is proportional to W and varies inversely as the square of the distance from the source.

The result we have just obtained applies to the heat flow in the vicinity of a point source of heat. Let's now try to find the equations that hold in the most general kind of heat flow, keeping only the condition that heat is conserved. We will be dealing only with what happens at places outside of any sources or absorbers of heat.

The differential equation for the conduction of heat was derived in Chapter 2. According to Eq. (2.44),

$$\mathbf{h} = -\kappa \nabla T. \quad (3.25)$$

(Remember that this relationship is an approximate one, but fairly good for some materials like metals.) It is applicable, of course, only in regions of the material where there is no generation or absorption of heat. We derived above another relation, Eq. (3.21), that holds when heat is conserved. If we combine that equation with (3.25), we get

$$-\frac{dq}{dt} = \nabla \cdot \mathbf{h} = -\nabla \cdot (\kappa \nabla T),$$

or

$$\frac{dq}{dt} = \kappa \nabla \cdot \nabla T = \kappa \nabla^2 T, \quad (3.26)$$

if κ is a constant. You remember that q is the amount of heat in a unit volume and $\nabla \cdot \nabla = \nabla^2$ is the Laplacian operator

$$\nabla^2 = \frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} + \frac{\partial^2}{\partial z^2}.$$

If we now make one more assumption we can obtain a very interesting equation. We assume that the temperature of the material is proportional to the heat content per unit volume—that is, that the material has a definite specific heat. When this assumption is valid (as it often is), we can write

$$\Delta q = c_v \Delta T$$

or

$$\frac{dq}{dt} = c_v \frac{dT}{dt}. \quad (3.27)$$

The rate of change of heat is proportional to the rate of change of temperature. The constant or proportionality c_v is, here, the specific heat per unit *volume* of the material. Using Eq. (3.27) with (3.26), we get

$$\frac{dT}{dt} = \frac{\kappa}{c_v} \nabla^2 T. \quad (3.28)$$

We find that the *time* rate of change of T —at every point—is proportional to the Laplacian of T , which is the second derivative of its spatial dependence. We have a differential equation—in x , y , z , and t —for the temperature T .

The differential equation (3.28) is called the *heat diffusion equation*. It is often written as

$$\frac{dT}{dt} = D \nabla^2 T, \quad (3.29)$$

where D is called the *diffusion constant*, and is here equal to κ/c_v .

The diffusion equation appears in many physical problems—in the diffusion of gases, in the diffusion of neutrons, and in others. We have already discussed the physics of some of these phenomena in Chapter 43 of Vol. I. Now you have the complete equation that describes diffusion in the most general possible situation. At some later time we will take up ways of solving the diffusion equation to find how the temperature varies in particular cases. We turn back now to consider other theorems about vector fields.

3-5 The circulation of a vector field

We wish now to look at the curl in somewhat the same way we looked at the divergence. We obtained Gauss' theorem by considering the integral over a surface, although it was not obvious at the beginning that we were going to be dealing with the divergence. How did we know that we were supposed to integrate over a surface in order to get the divergence? It was not at all clear that this would be the result. And so with an apparent equal lack of justification, we shall calculate something else about a vector and show that it is related to the curl. This time we calculate what is called the circulation of a vector field. If \mathbf{C} is any vector field, we take its component along a curved line and take the integral of this component all the way around a complete loop. The integral is called the *circulation* of the vector field around the loop. We have already considered a line integral of $\nabla\psi$ earlier in this chapter. Now we do the same kind of thing for *any* vector field \mathbf{C} .

Let Γ be any closed loop in space—imaginary, of course. An example is given in Fig. 3-7. The line integral of the tangential component of \mathbf{C} around the loop is written as

$$\oint_{\Gamma} C_t ds = \oint_{\Gamma} \mathbf{C} \cdot d\mathbf{s}. \quad (3.30)$$

You should note that the integral is taken all the way around, not from one point to another as we did before. The little circle on the integral sign is to remind us that the integral is to be taken all the way around. This integral is called the circulation of the vector field around the curve Γ . The name came originally from considering the circulation of a liquid. But the name—like flux—has been extended to apply to any field even when there is no material “circulating.”

Playing the same kind of game we did with the flux, we can show that the circulation around a loop is the sum of the circulations around two partial loops. Suppose we break up our curve of Fig. 3-7 into two loops, by joining two points (1) and (2) on the original curve by some line that cuts across as shown in Fig. 3-8. There are now two loops, Γ_1 and Γ_2 . Γ_1 is made up of Γ_a , which is that part of the original curve to the left of (1) and (2), plus Γ_{ab} , the “short cut.” Γ_2 is made up of the rest of the original curve plus the short cut.

The circulation around Γ_1 is the sum of an integral along Γ_a and along Γ_{ab} . Similarly, the circulation around Γ_2 is the sum of two parts, one along Γ_b and the other along Γ_{ab} . The integral along Γ_{ab} will have, for the curve Γ_2 , the opposite sign from what it has for Γ_1 , because the direction of travel is opposite—we must take both our line integrals with the same “sense” of rotation.

Following the same kind of argument we used before, you can see that the sum of the two circulations will give just the line integral around the original curve Γ . The parts due to Γ_{ab} cancel. The circulation around the one part plus the circulation around the second part equals the circulation about the outer line. We can continue the process of cutting the original loop into any number of smaller loops. When we add the circulations of the smaller loops, there is always a cancellation of the parts on their adjacent portions, so that the sum is equivalent to the circulation around the original single loop.

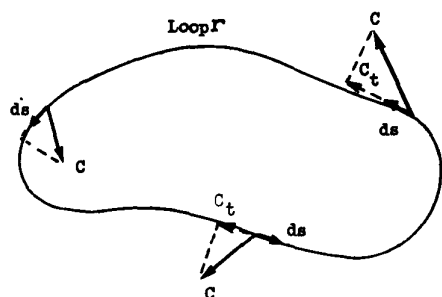


Fig. 3-7. The circulation of \mathbf{C} around the curve Γ is the line integral of C_t , the tangential component of \mathbf{C} .

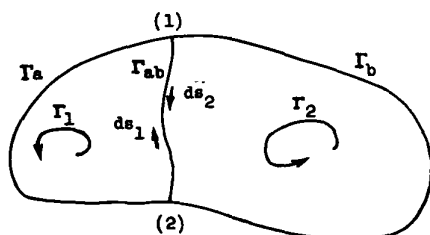


Fig. 3-8. The circulation around the whole loop is the sum of the circulations around the two loops: $\Gamma_1 = \Gamma_a + \Gamma_{ab}$ and $\Gamma_2 = \Gamma_b + \Gamma_{ab}$.

Now let us suppose that the original loop is the boundary of some surface. There are, of course, an infinite number of surfaces which all have the original loops as the boundary. Our results will not, however, depend on which surface we choose. First, we break our original loop into a number of small loops that all lie on the surface we have chosen, as in Fig. 3-9. No matter what the shape of the surface, if we choose our small loops small enough, we can assume that each of the small loops will enclose an area which is essentially flat. Also, we can choose our small loops so that each is very nearly a square. Now we can calculate the circulation around the big loop Γ by finding the circulations around all of the little squares and then taking their sum.

3-6 The circulation around a square; Stokes' theorem

How shall we find the circulation for each little square? One question is, how is the square oriented in space? We could easily make the calculation if it had a special orientation. For example, if it were in one of the coordinate planes. Since we have not assumed anything as yet about the orientation of the coordinate axes, we can just as well choose the axes so that the one little square we are concentrating on at the moment lies in the xy -plane, as in Fig. 3-10. If our result is expressed in vector notation, we can say that it will be the same no matter what the particular orientation of the plane.

We want now to find the circulation of the field C around our little square. It will be easy to do the line integral if we make the square small enough that the vector C doesn't change much along any one side of the square. (The assumption is better the smaller the square, so we are really talking about infinitesimal squares.) Starting at the point (x, y) —the lower left corner of the figure—we go around in the direction indicated by the arrows. Along the first side—marked (1)—the tangential component is $C_x(1)$ and the distance is Δx . The first part of the integral is $C_x(1) \Delta x$. Along the second leg, we get $C_y(2) \Delta y$. Along the third, we get $-C_x(3) \Delta x$, and along the fourth, $-C_y(4) \Delta y$. The minus signs are required because we want the tangential component in the direction of travel. The whole line integral is then

$$\oint C \cdot ds = C_x(1) \Delta x + C_y(2) \Delta y - C_x(3) \Delta x - C_y(4) \Delta y. \quad (3.31)$$

Now let's look at the first and third pieces. Together they are

$$[C_x(1) - C_x(3)] \Delta x. \quad (3.32)$$

You might think that to our approximation the difference is zero. That is true to the first approximation. We can be more accurate, however, and take into account the rate of change of C_x . If we do, we may write

$$C_x(3) = C_x(1) + \frac{\partial C_x}{\partial y} \Delta y. \quad (3.33)$$

If we included the next approximation, it would involve terms in $(\Delta y)^2$, but since we will ultimately think of the limit as $\Delta y \rightarrow 0$, such terms can be neglected. Putting (3.33) together with (3.32), we find that

$$[C_x(1) - C_x(3)] \Delta y = -\frac{\partial C_x}{\partial y} \Delta x \Delta y. \quad (3.34)$$

The derivative can, to our approximation, be evaluated at (x, y) .

Similarly, for the other two terms in the circulation, we may write

$$C_y(2) \Delta y - C_y(4) \Delta y = \frac{\partial C_y}{\partial x} \Delta x \Delta y. \quad (3.35)$$

The circulation around our square is then

$$\left(\frac{\partial C_y}{\partial x} - \frac{\partial C_x}{\partial y} \right) \Delta x \Delta y, \quad (3.36)$$

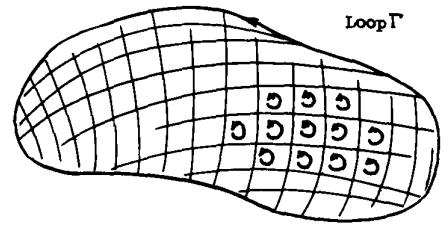


Fig. 3-9. Some surface bounded by the loop Γ is chosen. The surface is divided into a number of small areas, each approximately a square. The circulation around Γ is the sum of the circulations around the little loops.

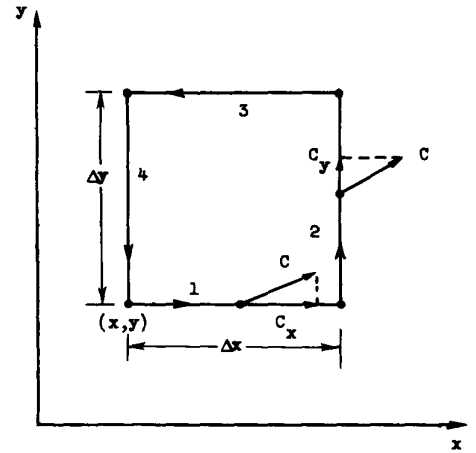


Fig. 3-10. Computing the circulation of C around a small square.

which is interesting, because the two terms in the parentheses are just the z -component of the curl. Also, we note that $\Delta x \Delta y$ is the area of our square. So we can write our circulation (3.36) as

$$(\nabla \times \mathbf{C})_z \Delta a.$$

But the z -component really means the component *normal* to the surface element. We can, therefore, write the circulation around a differential square in an invariant vector form:

$$\oint \mathbf{C} \cdot d\mathbf{s} = (\nabla \times \mathbf{C})_n \Delta a = (\nabla \times \mathbf{C}) \cdot \mathbf{n} \Delta a. \quad (3.37)$$

Our result is: the circulation of any vector \mathbf{C} around an infinitesimal square is the component of the curl of \mathbf{C} normal to the surface, times the area of the square.

The circulation around any loop Γ can now be easily related to the curl of the vector field. We fill in the loop with any convenient surface S , as in Fig. 3-11, and add the circulations around a set of infinitesimal squares in this surface. The sum can be written as an integral. Our result is a very useful theorem called Stokes' theorem (after Mr. Stokes).

STOKES' THEOREM.

$$\oint_{\Gamma} \mathbf{C} \cdot d\mathbf{s} = \int_S (\nabla \times \mathbf{C})_n da, \quad (3.38)$$

where S is any surface bounded by Γ .

We must now speak about a convention of signs. In Fig. 3-10 the z -axis would point *toward* you in a "usual"—that is, "right-handed"—system of axes. When we took our line integral with a "positive" sense of rotation, we found that the circulation was equal to the z -component of $\nabla \times \mathbf{C}$. If we had gone around the other way, we would have gotten the opposite sign. Now how shall we know, in general, what direction to choose for the positive direction of the "normal" component of $\nabla \times \mathbf{C}$? The "positive" normal must always be related to the sense of rotation, as in Fig. 3-10. It is indicated for the general case in Fig. 3-11.

One way of remembering the relationship is by the "right-hand rule." If you make the fingers of your *right* hand go around the curve Γ , with the fingertips pointed in the direction of the positive sense of $d\mathbf{s}$, then your thumb points in the direction of the *positive* normal to the surface S .

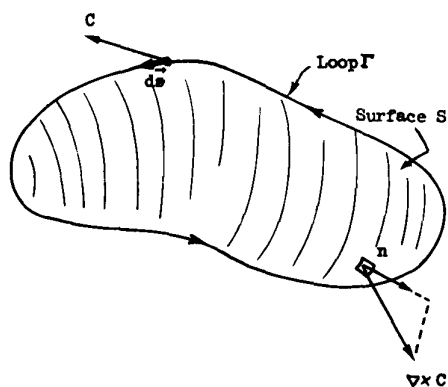


Fig. 3-11. The circulation of \mathbf{C} around Γ is the surface integral of the normal component of $\nabla \times \mathbf{C}$.

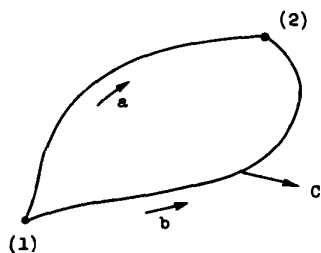


Fig. 3-12. If $\nabla \times \mathbf{C}$ is zero, the circulation around the closed curve Γ is zero. The line integral of $\mathbf{C} \cdot d\mathbf{s}$ from (1) to (2) along a must be the same as the line integral along b .

3-7 Curl-free and divergence-free fields

We would like, now, to consider some consequences of our new theorems. Take first the case of a vector whose curl is *everywhere* zero. Then Stokes' theorem says that the circulation around any loop is zero. Now if we choose two points (1) and (2) on a closed curve (Fig. 3-12), it follows that the line integral of the tangential component from (1) to (2) is independent of which of the two possible paths is taken. We can conclude that the integral from (1) to (2) can depend only on the location of these points—that is to say, it is some function of position only. The same logic was used in Chapter 14 of Vol. I, where we proved that if the integral around a closed loop of some quantity is always zero, then that integral can be represented as the difference of a function of the position of the two ends. This fact allowed us to invent the idea of a potential. We proved, furthermore, that the vector field was the gradient of this potential function (see Eq. 14.13 of Vol. I).

It follows that any vector field whose curl is zero is equal to the gradient of some scalar function. That is, if $\nabla \times \mathbf{C} = 0$, everywhere, there is some ψ (psi) for which $\mathbf{C} = \nabla\psi$ —a useful idea. We can, if we wish, describe this special kind of vector field by means of a scalar field.

Let's show something else. Suppose we have *any* scalar field ϕ (phi). If we take its gradient, $\nabla\phi$, the integral of this vector around any closed loop must be zero. Its line integral from point (1) to point (2) is $[\phi(2) - \phi(1)]$. If (1) and (2)

are the same points, our Theorem 1, Eq. (3.8), tells us that the line integral is zero:

$$\oint_{\text{loop}} \nabla \phi \cdot ds = 0.$$

Using Stokes' theorem, we can conclude that

$$\int \nabla \times (\nabla \phi) da = 0$$

over *any* surface. But if the integral is zero over *any* surface, the integrand must be zero. So

$$\nabla \times (\nabla \phi) = 0, \text{ always.}$$

We proved the same result in Section 2-7 by vector algebra.

Let's look now at a special case in which we fill in a *small* loop Γ with a *large* surface S , as indicated in Fig. 3-13. We would like, in fact, to see what happens when the loop shrinks down to a point, so that the surface boundary disappears—the surface becomes closed. Now if the vector C is everywhere finite, the line integral around Γ must go to zero as we shrink the loop—the integral is roughly proportional to the circumference of Γ , which goes to zero. According to Stokes' theorem, the surface integral of $(\nabla \times C)_n$ must also vanish. Somehow, as we close the surface we add in contributions that cancel out what was there before. So we have a new theorem:

$$\int_{\text{any closed surface}} (\nabla \times C)_n da = 0. \quad (3.39)$$

Now this is interesting, because we already have a theorem about the surface integral of a vector field. Such a surface integral is equal to the volume integral of the divergence of the vector, according to Gauss' theorem (Eq. 3.18). Gauss' theorem, applied to $\nabla \times C$, says

$$\int_{\text{closed surface}} (\nabla \times C)_n da = \int_{\text{volume inside}} \nabla \cdot (\nabla \times C) dV. \quad (3.40)$$

So we conclude that the second integral must also be zero:

$$\int_{\text{any volume}} \nabla \cdot (\nabla \times C) dV = 0, \quad (3.41)$$

and this is true for any vector field C whatever. Since Eq. (3.41) is true for *any* volume, it must be true that at *every* point in space the integrand is zero. We have

$$\nabla \cdot (\nabla \times C) = 0, \text{ always.}$$

But this is the same result we got from vector algebra in Section 2-7. Now we begin to see how everything fits together.

3-8 Summary

Let us summarize what we have found about the vector calculus. These are really the salient points of Chapters 2 and 3:

1. The operators $\partial/\partial x$, $\partial/\partial y$, and $\partial/\partial z$ can be considered as the three components of a vector operator ∇ , and the formulas which result from vector algebra by treating this operator as a vector are correct:

$$\nabla = \left(\frac{\partial}{\partial x}, \frac{\partial}{\partial y}, \frac{\partial}{\partial z} \right).$$

2. The difference of the values of a scalar field at two points is equal to the line integral of the tangential component of the gradient of that scalar along

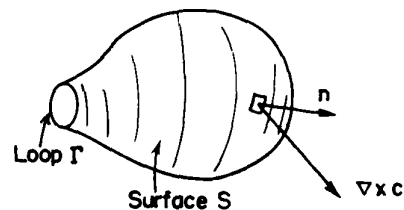


Fig. 3-13. Going to the limit of a closed surface, we find that the surface integral of $(\nabla \times C)_n$ must vanish.

any curve at all between the first and second points:

$$\psi(2) - \psi(1) = \int_{\substack{(1) \\ \text{any curve}}}^{(2)} \nabla \psi \cdot ds. \quad (3.42)$$

3. The surface integral of the normal component of an arbitrary vector over a closed surface is equal to the integral of the divergence of the vector over the volume interior to the surface:

$$\int_{\substack{\text{closed} \\ \text{surface}}} \mathbf{C} \cdot \mathbf{n} \, da = \int_{\substack{\text{volume} \\ \text{inside}}} \nabla \cdot \mathbf{C} \, dV. \quad (3.43)$$

4. The line integral of the tangential component of an arbitrary vector around a closed loop is equal to the surface integral of the normal component of the curl of that vector over any surface which is bounded by the loop.

$$\int_{\text{boundary}} \mathbf{C} \cdot ds = \int_{\text{surface}} (\nabla \times \mathbf{C}) \cdot \mathbf{n} \, da. \quad (3.44)$$

Electrostatics

4-1 Statics

We begin now our detailed study of the theory of electromagnetism. All of electromagnetism is contained in the Maxwell equations.

Maxwell's equations:

$$\nabla \cdot \mathbf{E} = \frac{\rho}{\epsilon_0}, \quad (4.1)$$

$$\nabla \times \mathbf{E} = -\frac{\partial \mathbf{B}}{\partial t}, \quad (4.2)$$

$$c^2 \nabla \times \mathbf{B} = \frac{\partial \mathbf{E}}{\partial t} + \frac{\mathbf{j}}{\epsilon_0}, \quad (4.3)$$

$$\nabla \cdot \mathbf{B} = 0. \quad (4.4)$$

The situations that are described by these equations can be very complicated. We will consider first relatively simple situations, and learn how to handle them before we take up more complicated ones. The easiest circumstance to treat is one in which nothing depends on the time—called the *static* case. All charges are permanently fixed in space, or if they do move, they move as a steady flow in a circuit (so ρ and \mathbf{j} are constant in time). In these circumstances, all of the terms in the Maxwell equations which are time derivatives of the field are zero. In this case, the Maxwell equations become:

Electrostatics:

$$\nabla \cdot \mathbf{E} = \frac{\rho}{\epsilon_0}, \quad (4.5)$$

$$\nabla \times \mathbf{E} = 0. \quad (4.6)$$

Magnetostatics:

$$\nabla \times \mathbf{B} = \frac{\mathbf{j}}{\epsilon_0 c^2}, \quad (4.7)$$

$$\nabla \cdot \mathbf{B} = 0. \quad (4.8)$$

You will notice an interesting thing about this set of four equations. It can be separated into two pairs. The electric field \mathbf{E} appears only in the first two, and the magnetic field \mathbf{B} appears only in the second two. The two fields are not interconnected. This means that *electricity and magnetism are distinct phenomena so long as charges and currents are static*. The interdependence of \mathbf{E} and \mathbf{B} does not appear until there are changes in charges or currents, as when a condenser is charged, or a magnet moved. Only when there are sufficiently rapid changes, so that the time derivatives in Maxwell's equations become significant, will \mathbf{E} and \mathbf{B} depend on each other.

Now if you look at the equations of statics you will see that the study of the two subjects we call electrostatics and magnetostatics is ideal from the point of view of learning about the mathematical properties of vector fields. Electrostatics is a neat example of a vector field with *zero curl* and a *given divergence*. Magnetostatics is a neat example of a field with *zero divergence* and a *given curl*. The more conventional—and you may be thinking, more satisfactory—way of presenting

4-1 Statics

4-2 Coulomb's law; superposition

4-3 Electric potential

4-4 $\mathbf{E} = -\nabla\phi$

4-5 The flux of \mathbf{E}

4-6 Gauss' law; the divergence of \mathbf{E}

4-7 Field of a sphere of charge

4-8 Field lines; equipotential surfaces

*Review: Chapters 13 and 14, Vol. I,
Work and Potential Energy*

$$\begin{aligned} \epsilon_0 c^2 &= \frac{10^7}{4\pi} \\ \frac{1}{4\pi\epsilon_0} &\approx 9 \times 10^9 \\ [\epsilon_0] &= \text{coulomb}^2/\text{newton-meter}^2 \end{aligned}$$

the theory of electromagnetism is to start first with electrostatics and thus to learn about the divergence. Magnetostatics and the curl are taken up later. Finally, electricity and magnetism are put together. We have chosen to start with the complete theory of vector calculus. Now we shall apply it to the special case of electrostatics, the field of E given by the first pair of equations.

We will begin with the simplest situations—ones in which the positions of all charges are specified. If we had only to study electrostatics at this level (as we shall do in the next two chapters), life would be very simple—in fact, almost trivial. Everything can be obtained from Coulomb's law and some integration, as you will see. In many real electrostatic problems, however, we do not *know*, initially, where the charges are. We know only that they have distributed themselves in ways that depend on the properties of matter. The positions that the charges take up depend on the E field, which in turn depends on the positions of the charges. Then things can get quite complicated. If, for instance, a charged body is brought near a conductor or insulator, the electrons and protons in the conductor or insulator will move around. The charge density ρ in Eq. (4.5) may have one part that we know about, from the charge that we brought up; but there will be other parts from charges that have moved around in the conductor. And all of the charges must be taken into account. One can get into some rather subtle and interesting problems. So although this chapter is to be on electrostatics, it will not cover the more beautiful and subtle parts of the subject. It will treat only the situation where we can assume that the positions of all the charges are known. Naturally, you should be able to do that case before you try to handle the other ones.

4-2 Coulomb's law; superposition

It would be logical to use Eqs. (4.5) and (4.6) as our starting points. It will be easier, however, if we start somewhere else and come back to these equations. The results will be equivalent. We will start with a law that we have talked about before, called Coulomb's law, which says that between two charges at rest there is a force directly proportional to the product of the charges and inversely proportional to the square of the distance between. The force is along the straight line from one charge to the other.

$$\text{Coulomb's law:} \quad F_1 = \frac{1}{4\pi\epsilon_0} \frac{q_1 q_2}{r_{12}^2} e_{12} = -F_2. \quad (4.9)$$

F_1 is the force *on* charge q_1 , e_{12} is the unit vector in the direction *to* q_1 *from* q_2 , and r_{12} is the distance between q_1 and q_2 . The force F_2 on q_2 is equal and opposite to F_1 .

The constant of proportionality, for historical reasons, is written as $1/4\pi\epsilon_0$. In the system of units which we use—the mks system—it is defined as exactly 10^{-7} times the speed of light squared. Now since the speed of light is approximately 3×10^8 meters per second, the constant is approximately 9×10^9 , and the unit turns out to be newton·meter²/coulomb² or volt·meter per coulomb.

$$\begin{aligned} \frac{1}{4\pi\epsilon_0} &= 10^{-7} c^2 \quad (\text{by definition}) \\ &= 9.0 \times 10^9 \quad (\text{by experiment}). \end{aligned} \quad (4.10)$$

Unit: newton·meter²/coulomb²,
or volt·meter/coulomb.

When there are more than two charges present—the only really interesting times—we must supplement Coulomb's law with one other fact of nature: the force on any charge is the vector sum of the Coulomb forces from each of the other charges. This fact is called “the principle of superposition.” That's all there is to electrostatics. If we combine the Coulomb law and the principle of superposition, there is nothing else. Equations (4.5) and (4.6)—the electrostatic equations—say no more and no less.

When applying Coulomb's law, it is convenient to introduce the idea of an electric field. We say that the field $E(1)$ is the force *per unit charge* on q_1 (due to all other charges). Dividing Eq. (4.9) by q_1 , we have, for one other charge besides q_1 ,

$$E(1) = \frac{1}{4\pi\epsilon_0} \frac{q_2}{r_{12}^2} e_{12}. \quad (4.11)$$

Also, we consider that $E(1)$ describes something about the point (1) even if q_1 were not there—assuming that all other charges keep their same positions. We say: $E(1)$ is the electric field *at* the point (1).

The electric field E is a vector, so by Eq. (4.11) we really mean three equations—one for each component. Writing out explicitly the x -component, Eq. (4.11) means

$$E_x(x_1, y_1, z_1) = \frac{q_2}{4\pi\epsilon_0} \frac{x_1 - x_2}{[(x_1 - x_2)^2 + (y_1 - y_2)^2 + (z_1 - z_2)^2]^{3/2}}, \quad (4.12)$$

and similarly for the other components.

If there are many charges present, the field E at any point (1) is a sum of the contributions from each of the other charges. Each term of the sum will look like (4.11) or (4.12). Letting q_j be the magnitude of the j th charge, and r_{1j} the displacement from q_j to the point (1), we write

$$E(1) = \sum_j \frac{1}{4\pi\epsilon_0} \frac{q_j}{r_{1j}^2} e_{1j}. \quad (4.13)$$

Which means, of course,

$$E_x(x_1, y_1, z_1) = \sum_j \frac{1}{4\pi\epsilon_0} \frac{q_j(x_1 - x_j)}{[(x_1 - x_j)^2 + (y_1 - y_j)^2 + (z_1 - z_j)^2]^{3/2}}, \quad (4.14)$$

and so on.

Often it is convenient to ignore the fact that charges come in packages like electrons and protons, and think of them as being spread out in a continuous smear—or in a “distribution,” as it is called. This is O.K. so long as we are not interested in what is happening on too small a scale. We describe a charge distribution by the “charge density,” $\rho(x, y, z)$. If the amount of charge in a small volume ΔV_2 located at the point (2) is Δq_2 , then ρ is defined by

$$\Delta q_2 = \rho(2) \Delta V_2. \quad (4.15)$$

To use Coulomb's law with such a description, we replace the sums of Eqs. (4.13) or (4.14) by integrals over all volumes containing charges. Then we have

$$E(1) = \frac{1}{4\pi\epsilon_0} \int_{\text{all space}} \frac{\rho(2) e_{12} dV_2}{r_{12}^2}. \quad (4.16)$$

Some people prefer to write

$$e_{12} = \frac{\mathbf{r}_{12}}{r_{12}},$$

where \mathbf{r}_{12} is the vector displacement *to* (1) *from* (2), as shown in Fig. 4-1. The integral for E is then written as

$$E(1) = \frac{1}{4\pi\epsilon_0} \int_{\text{all space}} \frac{\rho(2) \mathbf{r}_{12} dV_2}{r_{12}^3}. \quad (4.17)$$

When we want to calculate something with these integrals, we usually have to write them out in explicit detail. For the x -component of either Eq. (4.16) or (4.17), we would have

$$E_x(x_1, y_1, z_1) = \int_{\text{all space}} \frac{(x_1 - x_2) \rho(x_2, y_2, z_2) dx_2 dy_2 dz_2}{4\pi\epsilon_0 [(x_1 - x_2)^2 + (y_1 - y_2)^2 + (z_1 - z_2)^2]^{3/2}}. \quad (4.18)$$

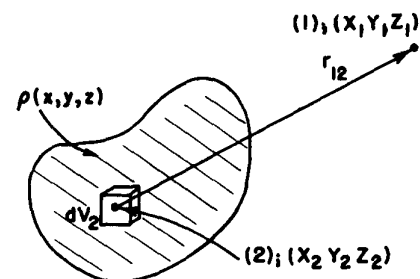


Fig. 4-1. The electric field E at point (1), from a charge distribution, is obtained from an integral over the distribution. Point (1) could also be inside the distribution.

We are not going to use this formula much. We write it here only to emphasize the fact that we have completely solved all the electrostatic problems in which we know the locations of all of the charges. Given the charges, what are the fields? *Answer:* Do this integral. So there is nothing to the subject; it is just a case of doing complicated integrals over three dimensions—strictly a job for a computing machine!

With our integrals we can find the fields produced by a sheet of charge, from a line of charge, from a spherical shell of charge, or from any specified distribution. It is important to realize, as we go on to draw field lines, to talk about potentials, or to calculate divergences, that we already have the answer here. It is merely a matter of it being sometimes easier to do an integral by some clever guesswork than by actually carrying it out. The guesswork requires learning all kinds of strange things. In practice, it might be easier to forget trying to be clever and always to do the integral directly instead of being so smart. We are, however, going to try to be smart about it. We shall go on to discuss some other features of the electric field.

4-3 Electric potential

First we take up the idea of electric potential, which is related to the work done in carrying a charge from one point to another. There is some distribution of charge, which produces an electric field. We ask about how much work it would take to carry a small charge from one place to another. The work done *against* the electrical forces in carrying a charge along some path is the *negative* of the component of the electrical force in the direction of the motion, integrated along the path. If we carry a charge from point *a* to point *b*,

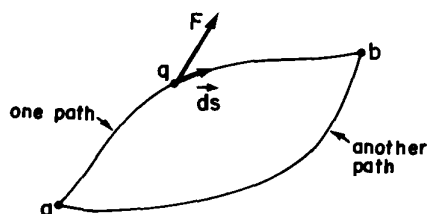


Fig. 4-2. The work done in carrying a charge from *a* to *b* is the negative of the integral of $F \cdot ds$ along the path taken.

$$W = - \int_a^b F \cdot ds,$$

where F is the electrical force *on* the charge at each point, and ds is the differential vector displacement along the path. (See Fig. 4-2.)

It is more interesting for our purposes to consider the work that would be done in carrying *one unit* of charge. Then the force on the charge is numerically the same as the electric field. Calling the work done against electrical forces in this case $W(\text{unit})$, we write

$$W(\text{unit}) = - \int_a^b E \cdot ds. \quad (4.19)$$

Now, in general, what we get with this kind of an integral depends on the path we take. But if the integral of (4.19) depended on the path from *a* to *b*, we could get work out of the field by carrying the charge to *b* along one path and then back to *a* on the other. We would go to *b* along the path for which W is smaller and *back* along the other, getting *out* more work than we put *in*.

There is nothing impossible, in principle, about getting energy out of a field. We shall, in fact, encounter fields where it is possible. It could be that as you move a charge you produce forces on the other part of the “machinery.” If the “machinery” moved against the force it would lose energy, thereby keeping the total energy in the world constant. For *electrostatics*, however, there is no such “machinery.” We know what the forces back on the sources of the field are. They are the Coulomb forces on the charges responsible for the field. If the other charges are fixed in position—as we assume in *electrostatics* only—these back forces can do no work on them. There is no way to get energy from them—provided, of course, that the principle of energy conservation works for electrostatic situations. We believe that it will work, but let’s just show that it must follow from Coulomb’s law of force.

We consider first what happens in the field due to a single charge q . Let point *a* be at the distance r_1 from q , and point *b* at r_2 . Now we carry a different charge, which we will call the “test” charge, and whose magnitude we choose to

be one unit, from a to b . Let's start with the easiest possible path to calculate. We carry our test charge first along the arc of a circle, then along a radius, as shown in part (a) of Fig. 4-3. Now on that particular path it is child's play to find the work done (otherwise we wouldn't have picked it). First, there is no work done at all on the path from a to a' . The field is radial (from Coulomb's law), so it is at right angles to the direction of motion. Next, on the path from a' to b , the field is in the direction of motion and varies as $1/r^2$. Thus the work done on the test charge in carrying it from a to b would be

$$-\int_a^b \mathbf{E} \cdot d\mathbf{s} = -\frac{q}{4\pi\epsilon_0} \int_{a'}^b \frac{dr}{r^2} = -\frac{q}{4\pi\epsilon_0} \left(\frac{1}{r_a} - \frac{1}{r_b} \right). \quad (4.20)$$

Now let's take another easy path. For instance, the one shown in part (b) of Fig. 4-3. It goes for awhile along an arc of a circle, then radially for awhile, then along an arc again, then radially, and so on. Every time we go along the circular parts, we do no work. Every time we go along the radial parts, we must just integrate $1/r^2$. Along the first radial stretch, we integrate from r_a to $r_{a'}$, then along the next radial stretch from $r_{a'}$ to $r_{a''}$, and so on. The sum of all these integrals is the same as a single integral directly from r_a to r_b . We get the same answer for this path that we did for the first path we tried. It is clear that we would get the same answer for *any* path which is made up of an arbitrary number of the same kinds of pieces.

What about smooth paths? Would we get the same answer? We discussed this point previously in Chapter 13 of Vol. I. Applying the same arguments used there, we can conclude that work done in carrying a unit charge from a to b is independent of the path.

$$W(\text{unit}) \Bigg\{ \begin{matrix} a \rightarrow b \\ \text{any} \\ \text{path} \end{matrix} \Bigg\} = - \int_a^b \mathbf{E} \cdot d\mathbf{s}.$$

Since the work done depends only on the endpoints, it can be represented as the difference between two numbers. We can see this in the following way. Let's choose a reference point P_0 and agree to evaluate our integral by using a path that always goes *by way of* point P_0 . Let $\phi(a)$ stand for the work done against the field in going *from* P_0 to point a , and let $\phi(b)$ be the work done in going *from* P_0 to point b (Fig. 4-4). The work in going *to* P_0 from a (on the way to b) is the negative of $\phi(a)$, so we have that

$$-\int_a^b \mathbf{E} \cdot d\mathbf{s} = \phi(b) - \phi(a). \quad (4.21)$$

Since only the difference in the function ϕ at two points is ever involved, we do not really have to specify the location of P_0 . Once we have chosen some reference point, however, a number ϕ is determined for *any* point in space; ϕ is then a *scalar field*. It is a function of x, y, z . We call this scalar function the *electrostatic potential* at any point.

Electrostatic potential:

$$\phi(P) = - \int_{P_0}^P \mathbf{E} \cdot d\mathbf{s}. \quad (4.22)$$

For convenience, we will often take the reference point at infinity. Then, for a single charge at the origin, the potential ϕ is given for any point (x, y, z) —using Eq. (4.20):

$$\phi(x, y, z) = \frac{q}{4\pi\epsilon_0} \frac{1}{r}. \quad (4.23)$$

The electric field from several charges can be written as the sum of the electric field from the first, from the second, from the third, etc. When we integrate the sum to find the potential we get a sum of integrals. Each of the integrals is the

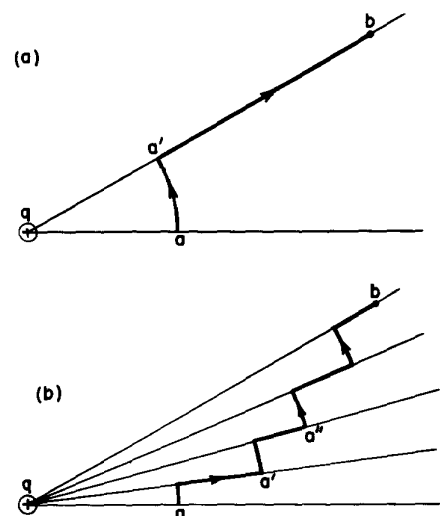


Fig. 4-3. In carrying a test charge from a to b the same work is done along either path.

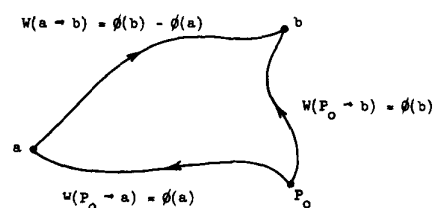


Fig. 4-4. The work done in going along any path from a to b is the negative of the work from some point P_0 to a plus the work from P_0 to b .

potential from one of the charges. We conclude that the potential ϕ from a lot of charges is the sum of the potentials from all the individual charges. There is a superposition principle also for potentials. Using the same kind of arguments by which we found the electric field from a group of charges and for a distribution of charges, we can get the complete formulas for the potential ϕ at a point we call (1):

$$\phi(1) = \sum_j \frac{1}{4\pi\epsilon_0} \frac{q_j}{r_{1j}}, \quad (4.24)$$

$$\phi(1) = \frac{1}{4\pi\epsilon_0} \int \frac{\rho(2) dV_2}{r_{12}}. \quad (4.25)$$

Remember that the potential ϕ has a physical significance: it is the potential energy which a unit charge would have if brought to the specified point in space from some reference point.

4-4 $E = -\nabla\phi$

Who cares about ϕ ? Forces on charges are given by E , the electric field. The point is that E can be obtained easily from ϕ —it is as easy, in fact, as taking a derivative. Consider two points, one at x and one at $(x + dx)$, but both at the same y and z , and ask how much work is done in carrying a unit charge from one point to the other. The path is along the horizontal line from x to $x + dx$. The work done is the difference in the potential at the two points:

$$\Delta W = \phi(x + \Delta x, y, z) - \phi(x, y, z) = \frac{\partial\phi}{\partial x} \Delta x.$$

But the work done against the field for the same path is

$$\Delta W = -\int E \cdot ds = -E_x \Delta x.$$

We see that

$$E_x = -\frac{\partial\phi}{\partial x}. \quad (4.26)$$

Similarly, $E_y = -\partial\phi/\partial y$, $E_z = -\partial\phi/\partial z$, or, summarizing with the notation of vector analysis,

$$E = -\nabla\phi. \quad (4.27)$$

This equation is the differential form of Eq. (4.22). Any problem with specified charges can be solved by computing the potential from (4.24) or (4.25) and using (4.27) to get the field. Equation (4.27) also agrees with what we found from vector calculus: that for any scalar field ϕ

$$\int_a^b \nabla\phi \cdot ds = \phi(b) - \phi(a). \quad (4.28)$$

According to Eq. (4.25) the scalar potential ϕ is given by a three-dimensional integral similar to the one we had for E . Is there any advantage to computing ϕ rather than E ? Yes. There is only one integral for ϕ , while there are three integrals for E —because it is a vector. Furthermore, $1/r$ is usually a little easier to integrate than x/r^3 . It turns out in many practical cases that it is easier to calculate ϕ and then take the gradient to find the electric field, than it is to evaluate the three integrals for E . It is merely a practical matter.

There is also a deeper physical significance to the potential ϕ . We have shown that E of Coulomb's law is obtained from $E = -\text{grad } \phi$, when ϕ is given by (4.22). But if E is equal to the gradient of a scalar field, then we know from the vector calculus that the curl of E must vanish:

$$\nabla \times E = 0. \quad (4.29)$$

But that is just our second fundamental equation of electrostatics, Eq. (4.6). We have shown that Coulomb's law gives an E field that satisfies that condition. So far, everything is all right.

We had really proved that $\nabla \times E$ was zero before we defined the potential. We had shown that the work done around a closed path is zero. That is, that

$$\oint E \cdot ds = 0$$

for *any* path. We saw in Chapter 3 that for any such field $\nabla \times E$ must be zero everywhere. The electric field in electrostatics is an example of a curl-free field.

You can practice your vector calculus by proving that $\nabla \times E$ is zero in a different way—by computing the components of $\nabla \times E$ for the field of a point charge, as given by Eq. (4.11). If you get zero, the superposition principle says you would get zero for the field of any charge distribution.

We should point out an important fact. For any *radial* force the work done is independent of the path, and there exists a potential. If you think about it, the entire argument we made above to show that the work integral was independent of the path depended only on the fact that the force from a single charge was radial and spherically symmetric. It did not depend on the fact that the dependence on distance was as $1/r^2$ —there could have been any r dependence. The existence of a potential, and the fact that the curl of E is zero, comes really only from the *symmetry* and *direction* of the electrostatic forces. Because of this, Eq. (4-28)—or (4.29)—can contain only part of the laws of electricity.

4-5 The flux of E

We will now derive a field equation that depends specifically and directly on the fact that the force law is inverse square. That the field varies inversely as the square of the distance seems, for some people, to be “only natural,” because “that’s the way things spread out.” Take a light source with light streaming out: the amount of light that passes through a surface cut out by a cone with its apex at the source is the same no matter at what radius the surface is placed. It must be so if there is to be conservation of light energy. The amount of light per unit area—the intensity—must vary inversely as the area cut by the cone, i.e., inversely as the square of the distance from the source. Certainly the electric field should vary inversely as the square of the distance for the same reason! But there is no such thing as the “same reason” here. Nobody can say that the electric field measures the flow of something like light which must be conserved. If we had a “model” of the electric field in which the electric field vector represented the direction and speed—say the current—of some kind of little “bullets” which were flying out, *and* if our model required that these bullets were conserved, that none could ever disappear once it was shot out of a charge, then we might say that we can “see” that the inverse square law is necessary. On the other hand, there would necessarily be some mathematical way to express this physical idea. If the electric field *were* like conserved bullets going out, then it would vary inversely as the square of the distance and we would be able to describe that behavior by an equation—which is purely mathematical. Now there is no harm in thinking this way, so long as we do not say that the electric field *is made* out of bullets, but realize that we are using a model to help us find the right mathematics.

Suppose, indeed, that we imagine for a moment that the electric field did represent the flow of something that was conserved—everywhere, that is, except at charges. (It has to start somewhere!) We imagine that whatever it is flows out of a charge into the space around. If E were the vector of such a flow (as h is for heat flow), it would have a $1/r^2$ dependence near a point source. Now we wish to use this model to find out how to state the inverse square law in a deeper or more abstract way, rather than simply saying “inverse square.” (You may wonder why we should want to avoid the direct statement of such a simple law, and want instead to imply the same thing sneakily in a different way. Patience! It will turn out to be useful.)

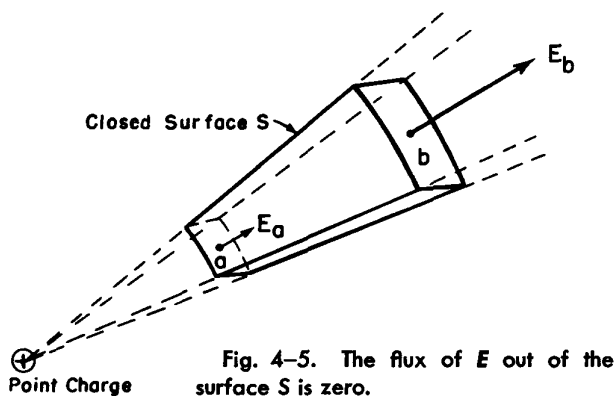


Fig. 4-5. The flux of E out of the surface S is zero.

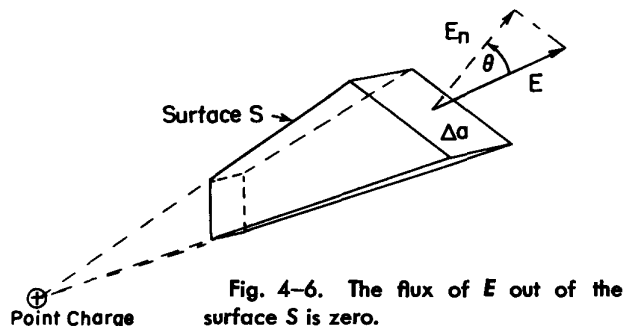


Fig. 4-6. The flux of E out of the surface S is zero.

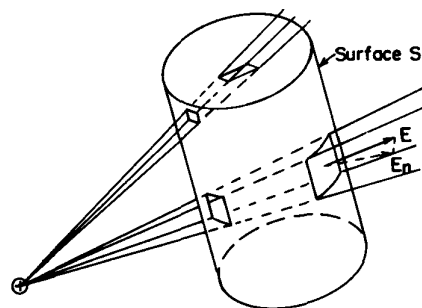


Fig. 4-7. Any volume can be thought of as completely made up of infinitesimal truncated cones. The flux of E from one end of each conical segment is equal and opposite to the flux from the other end. The total flux from the surface S is therefore zero.

We ask: What is the “flow” of E out of an arbitrary closed surface in the neighborhood of a point charge? First let’s take an easy surface—the one shown in Fig. 4-5. If the E field is like a flow, the net flow out of this box should be zero. That is what we get if by the “flow” from this surface we mean the surface integral of the normal component of E —that is, the flux of E . On the radial faces, the normal component is zero. On the spherical faces, the normal component E_n is just the magnitude of E —minus for the smaller face and plus for the larger face. The magnitude of E decreases as $1/r^2$, but the surface area is proportional to r^2 , so the product is independent of r . The flux of E into face a is just cancelled by the flux out of face b . The total flow out of S is zero, which is to say that for this surface

$$\int_S E_n da = 0. \quad (4.30)$$

Next we show that the two end surfaces may be tilted with respect to the radial line without changing the integral (4.30). Although it is true in general, for our purposes it is only necessary to show that this is true when the end surfaces are small, so that they subtend a small angle from the source—in fact, an infinitesimal angle. In Fig. 4-6 we show a surface S whose “sides” are radial, but whose “ends” are tilted. The end surfaces are not small in the figure, but you are to imagine the situation for very small end surfaces. Then the field E will be sufficiently uniform over the surface that we can use just its value at the center. When we tilt the surface by an angle θ , the area is increased by the factor $1/\cos \theta$. But E_n , the component of E normal to the surface, is decreased by the factor $\cos \theta$. The product $E_n \Delta a$ is unchanged. The flux out of the whole surface S is still zero.

Now it is easy to see that the flux out of a volume enclosed by *any* surface S must be zero. Any volume can be thought of as made up of pieces, like that in Fig. 4-6. The surface will be subdivided completely into pairs of end surfaces, and since the fluxes in and out of these end surfaces cancel by pairs, the total flux out of the surface will be zero. The idea is illustrated in Fig. 4-7. We have the completely general result that the total flux of E out of *any* surface S in the field of a point charge is zero.

But notice! Our proof works only if the surface S does not surround the charge. What would happen if the point charge were *inside* the surface? We could still divide our surface into pairs of areas that are matched by radial lines through the charge, as shown in Fig. 4-8. The fluxes through the two surfaces are still equal—by the same arguments as before—only now they have the *same* sign. The flux out of a surface that *surrounds* a charge is *not* zero. Then what is it? We can find out by a little trick. Suppose we “remove” the charge from the “inside” by surrounding the charge by a little surface S' totally inside the original surface S , as shown in Fig. 4-9. Now the volume enclosed *between* the two surfaces S and S' has no charge in it. The total flux out of this volume (including that through S') is zero, by the arguments we have given above. The arguments tell us, in fact, that the flux *into* the volume through S' is the same as the flux outward through S .

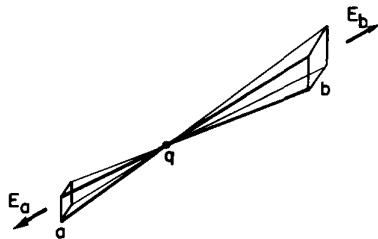


Fig. 4-8. If a charge is inside a surface, the flux out is not zero.

We can choose any shape we wish for S' , so let's make it a sphere centered on the charge, as in Fig. 4-10. Then we can easily calculate the flux through it. If the radius of the little sphere is r , the value of E everywhere on its surface is

$$\frac{1}{4\pi\epsilon_0} \frac{q}{r^2},$$

and is directed always normal to the surface. We find the total flux through S' if we multiply this normal component of E by the surface area:

$$\text{Flux through the surface } S' = \left(\frac{1}{4\pi\epsilon_0} \frac{q}{r^2} \right) (4\pi r^2) = \frac{q}{\epsilon_0}, \quad (4.31)$$

a number independent of the radius of the sphere! We know then that the flux outward through S is also q/ϵ_0 —a value independent of the shape of S so long as the charge q is inside.

We can write our conclusions as follows:

$$\int_{\text{any surface } S} E_n da = \begin{cases} 0; & q \text{ outside } S \\ \frac{q}{\epsilon_0}; & q \text{ inside } S \end{cases} \quad (4.32)$$

Let's return to our "bullet" analogy and see if it makes sense. Our theorem says that the net flow of bullets through a surface is zero if the surface does not enclose the gun that shoots the bullets. If the gun is enclosed in a surface, whatever size and shape it is, the number of bullets passing through is the same—it is given by the rate at which bullets are generated at the gun. It all seems quite reasonable for conserved bullets. But does the model tell us anything more than we get simply by writing Eq. (4.32)? No one has succeeded in making these "bullets" do anything else but produce this one law. After that, they produce nothing but errors. That is why today we prefer to represent the electromagnetic field purely abstractly.

4-6 Gauss' law; the divergence of E

Our nice result, Eq. (4.32), was proved for a single point charge. Now suppose that there are two charges, a charge q_1 at one point and a charge q_2 at another. The problem looks more difficult. The electric field whose normal component we integrate for the flux is the field due to both charges. That is, if E_1 represents the electric field that would have been produced by q_1 alone, and E_2 represents the electric field produced by q_2 alone, the total electric field is $E = E_1 + E_2$. The flux through any closed surface S is

$$\int_S (E_{1n} + E_{2n}) da = \int_S E_{1n} da + \int_S E_{2n} da. \quad (4.33)$$

The flux with both charges present is the flux due to a single charge plus the flux due to the other charge. If both charges are outside S , the flux through S is zero. If q_1 is inside S but q_2 is outside, then the first integral gives q_1/ϵ_0 and the second integral gives zero. If the surface encloses both charges, each will give its contribution and we have that the flux is $(q_1 + q_2)/\epsilon_0$. The general rule is clearly that the total flux out of a closed surface is equal to the total charge *inside*, divided by ϵ_0 .

Our result is an important general law of the electrostatic field, called Gauss' law.

$$\text{Gauss' law: } \int_{\text{any closed surface } S} E_n da = \frac{\text{sum of charges inside}}{\epsilon_0}, \quad (4.34)$$

$$\text{or } \int_{\text{any closed surface } S} \mathbf{E} \cdot \mathbf{n} da = \frac{Q_{\text{int}}}{\epsilon_0}, \quad (4.35)$$

where

$$Q_{\text{int}} = \sum_{\text{inside } S} q_i. \quad (4.36)$$

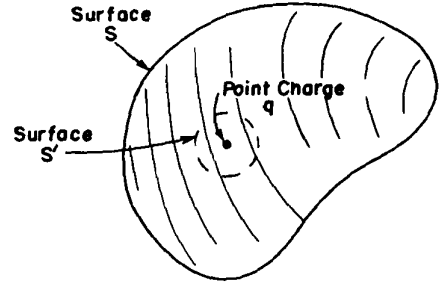


Fig. 4-9. The flux through S is the same as the flux through S' .

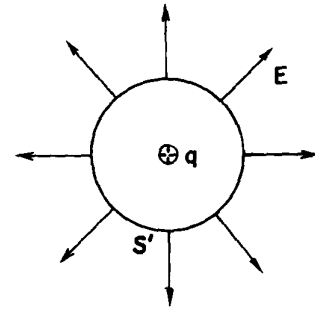


Fig. 4-10. The flux through a spherical surface containing a point charge q is q/ϵ_0 .

If we describe the location of charges in terms of a charge density ρ , we can consider that each infinitesimal volume dV contains a “point” charge ρdV . The sum over all charges is then the integral

$$Q_{\text{int}} = \int_{\text{volume inside } S} \rho dV. \quad (4.37)$$

From our derivation you see that Gauss’ law follows from the fact that the exponent in Coulomb’s law is exactly two. A $1/r^3$ field, or any $1/r^n$ field with $n \neq 2$, would not give Gauss’ law. So Gauss’ law is just an expression, in a different form, of the Coulomb law of forces between two charges. In fact, working back from Gauss’ law, you can derive Coulomb’s law. The two are quite equivalent so long as we keep in mind the rule that the forces between charges is radial.

We would now like to write Gauss’ law in terms of derivatives. To do this, we apply Gauss’ law to an infinitesimal cubical surface. We showed in Chapter 3 that the flux of E out of such a cube is $\nabla \cdot E$ times the volume dV of the cube. The charge inside of dV , by the definition of ρ , is equal to ρdV , so Gauss’ law gives

$$\nabla \cdot E dV = \frac{\rho dV}{\epsilon_0},$$

or

$$\nabla \cdot E = \frac{\rho}{\epsilon_0}. \quad (4.38)$$

The differential form of Gauss’ law is the first of our fundamental field equations of electrostatics, Eq. (4.5). We have now shown that the two equations of electrostatics, Eqs. (4.5) and (4.6), are equivalent to Coulomb’s law of force. We will now consider one example of the use of Gauss’ law. (We will come later to many more examples.)

4-7 Field of a sphere of charge

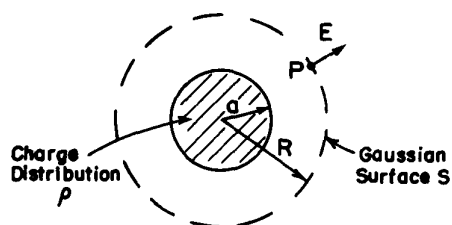


Fig. 4-11. Using Gauss’ law to find the field of a uniform sphere of charge.

One of the difficult problems we had when we studied the theory of gravitational attractions was to prove that the force produced by a solid sphere of matter was the same at the surface of the sphere as it would be if all the matter were concentrated at the center. For many years Newton didn’t make public his theory of gravitation, because he couldn’t be sure this theorem was true. We proved the theorem in Chapter 13 of Vol. I by doing the integral for the potential and then finding the gravitational force by using the gradient. Now we can prove the theorem in a most simple fashion. Only this time we will prove the corresponding theorem for a uniform sphere of electrical charge. (Since the laws of electrostatics are the same as those of gravitation, the same proof could be done for the gravitational field.)

We ask: What is the electric field E at a point P anywhere outside the surface of a sphere filled with a uniform distribution of charge? Since there is no “special” direction, we can assume that E is everywhere directed away from the center of the sphere. We consider an imaginary surface that is spherical and concentric with the sphere of charge, and that passes through the point P (Fig. 4-11). For this surface, the flux outward is

$$\int E_n da = E \cdot 4\pi R^2.$$

Gauss’ law tells us that this flux is equal to the total charge Q of the sphere (over ϵ_0):

$$E \cdot 4\pi R^2 = \frac{Q}{\epsilon_0},$$

or

$$E = \frac{1}{4\pi\epsilon_0} \frac{Q}{r^2}, \quad (4.39)$$

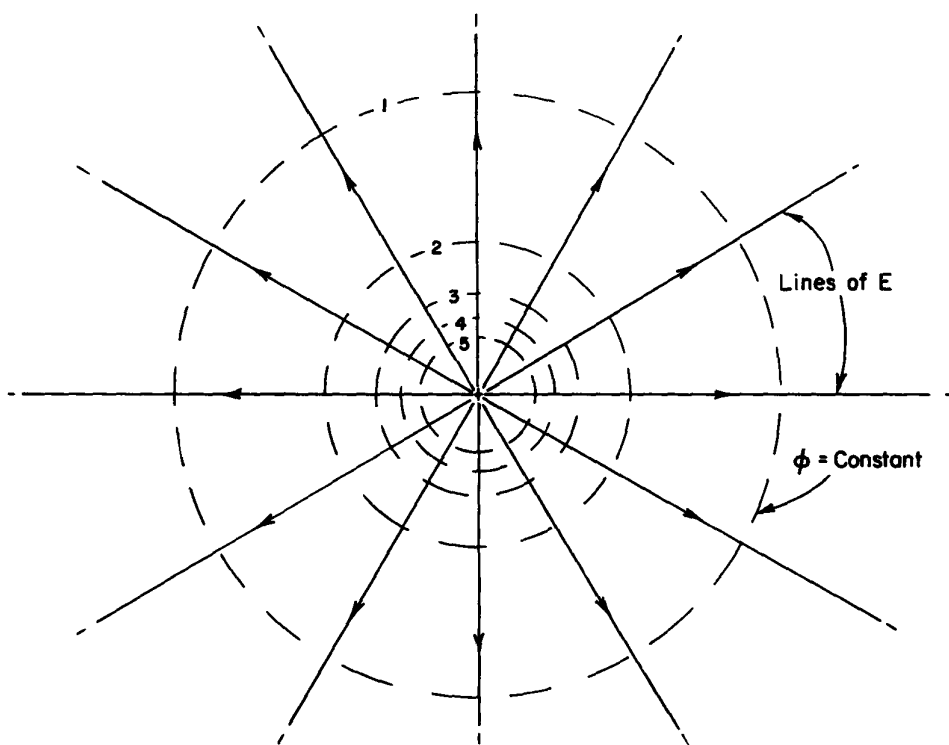


Fig. 4-12. Field lines and equipotential surfaces for a positive point charge.

which is the same formula we would have for a point charge Q . We have proved Newton's problem more easily than by doing the integral. It is, of course, a false kind of easiness—it has taken you some time to be able to understand Gauss' law, so you may think that no time has really been saved. But after you have used the theorem more and more, it begins to pay. It is a question of efficiency.

4-8 Field lines; equipotential surfaces

We would like now to give a geometrical description of the electrostatic field. The two laws of electrostatics, one that the flux is proportional to the charge inside and the other that the electric field is the gradient of a potential, can also be represented geometrically. We illustrate this with two examples.

First, we take the field of a point charge. We draw lines in the direction of the field—lines which are always tangent to the field, as in Fig. 4-12. These are called *field lines*. The lines show everywhere the direction of the electric vector. But we also wish to represent the magnitude of the vector. We can make the rule that the strength of the electric field will be represented by the "density" of the lines. By the density of the lines we mean the number of lines per unit area through a surface perpendicular to the lines. With these two rules we can have a picture of the electric field. For a point charge, the density of the lines must decrease as $1/r^2$. But the area of a spherical surface perpendicular to the lines at any radius r increases as r^2 , so if we always keep the same *number* of lines for *all* distances from the charge, the *density* will remain in proportion to the magnitude of the field. We can guarantee that there are the same number of lines at every distance if we insist that the lines be *continuous*—that once a line is started from the charge, it never stops. In terms of the field lines, Gauss' law says that lines should start only at plus charges and stop at minus charges. The number which *leave* a charge q must be equal to q/ϵ_0 .

Now, we can find a similar geometrical picture for the potential ϕ . The easiest way to represent the potential is to draw surfaces on which ϕ is a constant. We call them *equipotential* surfaces—surfaces of equal potential. Now what is the geometri-

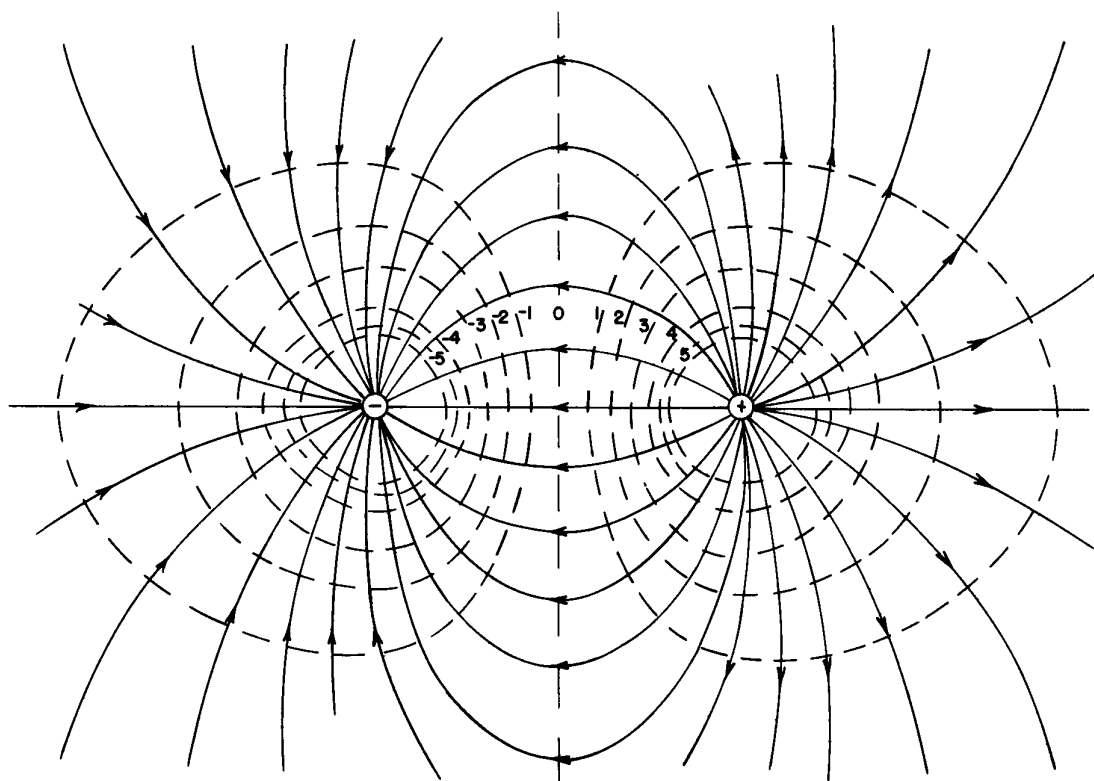


Fig. 4-13. Field lines and equipotentials for two equal and opposite point charges.

cal relationship of the equipotential surfaces to the field lines? The electric field is the gradient of the potential. The gradient is in the direction of the most rapid change of the potential, and is therefore perpendicular to an equipotential surface. If E were *not* perpendicular to the surface, it would have a component *in* the surface. The potential would be changing in the surface, but then it wouldn't be an equipotential. The equipotential surfaces must then be everywhere at right angles to the electric field lines.

For a point charge all by itself, the equipotential surfaces are spheres centered at the charge. We have shown in Fig. 4-12 the intersection of these spheres with a plane through the charge.

As a second example, we consider the field near two equal charges, a positive one and a negative one. To get the field is easy. The field is the superposition of the fields from each of the two charges. So, we can take two pictures like Fig. 4-12 and superimpose them—impossible! Then we would have field lines crossing each other, and that's not possible, because E can't have *two* directions at the same point. The disadvantage of the field-line picture is now evident. By geometrical arguments it is impossible to analyze in a very simple way where the new lines go. From the two independent pictures, we can't get the combined picture. The principle of superposition, a simple and deep principle about electric fields, does not have, in the field-line picture, an easy representation.

The field-line picture has its uses, however, so we might still like to draw the picture for a pair of equal (and opposite) charges. If we calculate the fields from Eq. (4.13) and the potentials from (4.23), we can draw the field lines and equipotentials. Figure 4-13 shows the result. But we first had to solve the problem mathematically!

A Note about Units

Quantity	Unit
F	newton
Q	coulomb
L	meter
W	joule
$\rho \sim Q/L^3$	coulomb/meter ³
$1/\epsilon_0 \sim FL^2/Q^2$	newton·meter ² /coulomb ²
$E \sim F/Q$	newton/coulomb
$\phi \sim W/Q$	joule/coulomb = volt
$E \sim \phi/L$	volt/meter
$1/\epsilon_0 \sim EL^2/Q$	volt·meter/coulomb

Application of Gauss' Law

5-1 Electrostatics is Gauss' law plus . . .

There are two laws of electrostatics: that the flux of the electric field from a volume is proportional to the charge inside—Gauss' law, and that the circulation of the electric field is zero— E is a gradient. From these two laws, all the predictions of electrostatics follow. But to say these things mathematically is one thing; to use them easily, and with a certain amount of ingenuity, is another. In this chapter we will work through a number of calculations which can be made with Gauss' law directly. We will prove theorems and describe some effects, particularly in conductors, that can be understood very easily from Gauss' law. Gauss' law by itself cannot give the solution of any problem because the other law must be obeyed too. So when we use Gauss' law for the solution of particular problems, we will have to add something to it. We will have to presuppose, for instance, some idea of how the field looks—based, for example, on arguments of symmetry. Or we may have to introduce specifically the idea that the field is the gradient of a potential.

5-2 Equilibrium in an electrostatic field

Consider first the following question: When can a point charge be in stable mechanical equilibrium in the electric field of other charges? As an example, imagine three negative charges at the corners of an equilateral triangle in a horizontal plane. Would a positive charge placed at the center of the triangle remain there? (It will be simpler if we ignore gravity for the moment, although including it would not change the results.) The force on the positive charge is zero, but is the equilibrium stable? Would the charge return to the equilibrium position if displaced slightly? The answer is no.

There are *no* points of stable equilibrium in *any* electrostatic field—except right on top of another charge. Using Gauss' law, it is easy to see why. First, for a charge to be in equilibrium at any particular point P_0 , the field must be zero. Second, if the equilibrium is to be a stable one, we require that if we move the charge away from P_0 in *any* direction, there should be a restoring force directed opposite to the displacement. The electric field at *all* nearby points must be pointing inward—toward the point P_0 . But that is in violation of Gauss' law if there is no charge at P_0 , as we can easily see.

Consider a tiny imaginary surface that encloses P_0 , as in Fig. 5-1. If the electric field everywhere in the vicinity is pointed toward P_0 , the surface integral of the normal component is certainly not zero. For the case shown in the figure, the flux through the surface must be a negative number. But Gauss' law says that the flux of electric field through any surface is proportional to the total charge inside. If there is no charge at P_0 , the field we have imagined violates Gauss' law. It is impossible to balance a positive charge in empty space—at a point where there is not some negative charge. A positive charge *can* be in equilibrium if it is in the middle of a distributed negative charge. Of course, the negative charge distribution would have to be held in place by other than electrical forces!

Our result has been obtained for a point charge. Does the same conclusion hold for a complicated arrangement of charges held together in fixed relative positions—with rods, for example? We consider the question for two equal charges fixed on a rod. Is it possible that this combination can be in equilibrium in some electrostatic field? The answer is again no. The *total* force on the rod cannot be restoring for displacements in every direction.

- 5-1 Electrostatics is Gauss' law plus . . .
- 5-2 Equilibrium in an electrostatic field
- 5-3 Equilibrium with conductors
- 5-4 Stability of atoms
- 5-5 The field of a line charge
- 5-6 A sheet of charge; two sheets
- 5-7 A sphere of charge; a spherical shell
- 5-8 Is the field of a point charge exactly $1/r^2$?
- 5-9 The fields of a conductor
- 5-10 The field in a cavity of a conductor

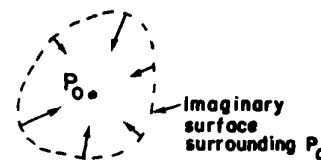


Fig. 5-1. If P_0 were a position of stable equilibrium for a positive charge, the electric field everywhere in the neighborhood would point toward P_0 .

Call F the total force on the rod in any position— F is then a vector field. Following the argument used above, we conclude that at a position of stable equilibrium, the divergence of F must be a negative number. But the total force on the rod is the first charge times the field at its position, plus the second charge times the field at its position:

$$F = q_1 E_1 + q_2 E_2. \quad (5.1)$$

The divergence of F is given by

$$\nabla \cdot F = q_1 (\nabla \cdot E_1) + q_2 (\nabla \cdot E_2).$$

If each of the two charges q_1 and q_2 is in free space, both $\nabla \cdot E_1$ and $\nabla \cdot E_2$ are zero, and $\nabla \cdot F$ is zero—not negative, as would be required for equilibrium. You can see that an extension of the argument shows that no rigid combination of any number of charges can have a position of stable equilibrium in an electrostatic field in free space.

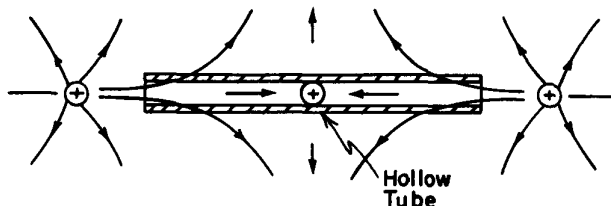


Fig. 5-2. A charge can be in equilibrium if there are mechanical constraints.

Now we have not shown that equilibrium is forbidden if there are pivots or other mechanical constraints. As an example, consider a hollow tube in which a charge can move back and forth freely, but not sideways. Now it is very easy to devise an electric field that points inward at both ends of the tube if it is allowed that the field may point laterally outward near the center of the tube. We simply place positive charges at each end of the tube, as in Fig. 5-2. There can now be an equilibrium point even though the divergence of E is zero. The charge, of course, would not be in stable equilibrium for sideways motion were it not for “non-electrical” forces from the tube walls.

5-3 Equilibrium with conductors

There is no stable spot in the field of a system of fixed charges. What about a system of charged conductors? Can a system of charged conductors produce a field that will have a stable equilibrium point for a point charge? (We mean at a point other than on a conductor, of course.) You know that conductors have the property that charges can move freely around in them. Perhaps when the point charge is displaced slightly, the other charges on the conductors will move in a way that will give a restoring force to the point charge? The answer is still no—although the proof we have just given doesn’t show it. The proof for this case is more difficult, and we will only indicate how it goes.

First, we note that when charges redistribute themselves on the conductors, they can only do so if their motion decreases their total potential energy. (Some energy is lost to heat as they move in the conductor.) Now we have already shown that if the charges producing a field are *stationary*, there is, near any zero point P_0 in the field, some direction for which moving a point charge away from P_0 will *decrease* the energy of the system (since the force is *away* from P_0). Any readjustment of the charges on the conductors can only lower the potential energy still more, so (by the principle of virtual work) their motion will only *increase* the force in that particular direction away from P_0 , and not reverse it.

Our conclusions do not mean that it is not possible to balance a charge by electrical forces. It is possible if one is willing to control the locations or the sizes of the supporting charges with suitable devices. You know that a rod standing on its point in a gravitational field is unstable, but this does not prove that it cannot be balanced on the end of a finger. Similarly, a charge can be held in one spot by electric fields if they are *variable*. But not with a passive—that is, a *static*—system.

5-4 Stability of atoms

If charges cannot be held stably in position, it is surely not proper to imagine matter to be made up of static *point* charges (electrons and protons) governed only by the laws of electrostatics. Such a static configuration is impossible; it would collapse!

It was once suggested that the positive charge of an atom could be distributed uniformly in a sphere, and the negative charges, the electrons, could be at rest inside the positive charge, as shown in Fig. 5-3. This was the first atomic model, proposed by Thompson. But Rutherford concluded from the experiment of Geiger and Marsden that the positive charges were very much concentrated, in what he called the nucleus. Thompson's static model had to be abandoned. Rutherford and Bohr then suggested that the equilibrium might be dynamic, with the electrons revolving in orbits, as shown in Fig. 5-4. The electrons would be kept from falling in toward the nucleus by their orbital motion. We already know at least one difficulty with this picture. With such motion, the electrons would be accelerating (because of the circular motion) and would, therefore, be radiating energy. They would lose the kinetic energy required to stay in orbit, and would spiral in toward the nucleus. Again unstable!

The stability of the atoms is now explained in terms of quantum mechanics. The electrostatic forces pull the electron as close to the nucleus as possible, but the electron is compelled to stay spread out in space over a distance given by the uncertainty principle. If it were confined in too small a space, it would have a great uncertainty in momentum. But that means that it would have a high expected energy—which it would use to escape from the electrical attraction. The net result is an electrical equilibrium not too different from the idea of Thompson—only it is the *negative* charge that is spread out (because the mass of the electron is so much smaller than the mass of the proton).

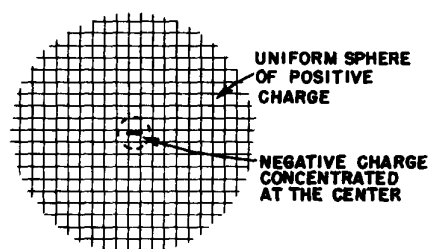


Fig. 5-3. The Thompson model of an atom.

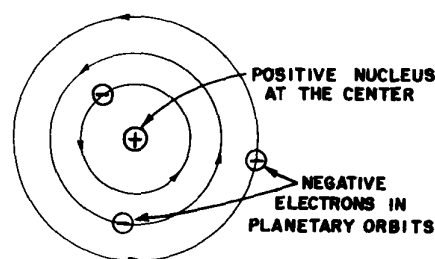


Fig. 5-4. The Rutherford-Bohr model of an atom.

5-5 The field of a line charge

Gauss' law can be used to solve a number of electrostatic field problems involving a special symmetry—usually spherical, cylindrical, or planar symmetry. In the remainder of this chapter we will apply Gauss' law to a few such problems. The ease with which these problems can be solved may give the misleading impression that the method is very powerful, and that one should be able to go on to many other problems. It is unfortunately not so. One soon exhausts the list of problems that can be solved easily with Gauss' law. In later chapters we will develop more powerful methods for investigating electrostatic fields.

As our first example, we consider a system with cylindrical symmetry. Suppose that we have a very long, uniformly charged rod. By this we mean that electric charges are distributed uniformly along an indefinitely long straight line, with the charge λ per unit length. We wish to know the electric field. The problem can, of course, be solved by integrating the contribution to the field from every part of the line. We are going to do it without integrating, by using Gauss' law and some guesswork. First, we surmise that the electric field will be directed radially outward from the line. Any axial component from charges on one side would be accompanied by an equal axial component from charges on the other side. The result could only be a radial field. It also seems reasonable that the field should have the same magnitude at all points equidistant from the line. This is obvious. (It may not be easy to prove, but it is true if space is symmetric—as we believe it is.)

We can use Gauss' law in the following way. We consider an *imaginary* surface in the shape of a cylinder coaxial with the line, as shown in Fig. 5-5. According to Gauss' law, the total flux of E from this surface is equal to the charge inside divided by ϵ_0 . Since the field is assumed to be normal to the surface, the normal component is the magnitude of the field. Let's call it E . Also, let the radius of the cylinder be r , and its length be taken as one unit, for convenience. The flux through the cylindrical surface is equal to E times the area of the surface, which is $2\pi r$. The flux through the two end faces is zero because the electric field is tan-

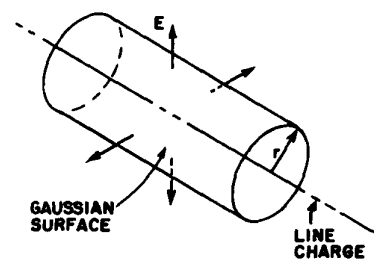


Fig. 5-5. A cylindrical gaussian surface coaxial with a line charge.

gential to them. The total charge inside our surface is just λ , because the length of the line inside is one unit. Gauss' law then gives

$$E \cdot 2\pi r = \lambda/\epsilon_0,$$

$$E = \frac{\lambda}{2\pi\epsilon_0 r}. \quad (5.2)$$

The electric field of a line charge depends inversely on the *first* power of the distance from the line.

5-6 A sheet of charge; two sheets

As another example, we will calculate the field from a uniform plane sheet of charge. Suppose that the sheet is infinite in extent and that the charge per unit area is σ . We are going to take another guess. Considerations of symmetry lead us to believe that the field direction is everywhere normal to the plane, and *if we have no field from any other charges in the world*, the fields must be the same (in magnitude) on each side. This time we choose for our Gaussian surface a rectangular box that cuts through the sheet, as shown in Fig. 5-6. The two faces parallel to the sheet will have equal areas, say A . The field is normal to these two faces, and parallel to the other four. The total flux is E times the area of the first face, plus E times the area of the opposite face—with no contribution from the other four faces. The total charge enclosed in the box is σA . Equating the flux to the charge inside, we have

$$EA + EA = \frac{\sigma A}{\epsilon_0},$$

$$E = \frac{\sigma}{2\epsilon_0}, \quad (5.3)$$

from which

a simple but important result.

You may remember that the same result was obtained in an earlier chapter by an integration over the entire surface. Gauss' law gives us the answer, in this instance, much more quickly (although it is not as generally applicable as the earlier method).

We emphasize that this result applies *only* to the field due to the charges on the sheet. If there are other charges in the neighborhood, the total field near the sheet would be the sum of (5.3) and the field of the other charges. Gauss' law would then tell us only that

$$E_1 + E_2 = \frac{\sigma}{\epsilon_0}, \quad (5.4)$$

where E_1 and E_2 are the fields directed outward on each side of the sheet.

The problem of two parallel sheets with equal and opposite charge densities, $+\sigma$ and $-\sigma$, is equally simple if we assume again that the outside world is quite symmetric. Either by superposing two solutions for a single sheet or by constructing a gaussian box that includes both sheets, it is easily seen that the field is zero *outside* of the two sheets (Fig. 5-7a). By considering a box that includes only one surface or the other, as in (b) or (c) of the figure, it can be seen that the field between the sheets must be twice what it is for a single sheet. The result is

$$E \text{ (between the sheets)} = \sigma/\epsilon_0, \quad (5.5)$$

$$E \text{ (outside)} = 0. \quad (5.6)$$

5-7 A sphere of charge; a spherical shell

We have already (in Chapter 4) used Gauss' law to find the field outside a uniformly charged spherical region. The same method can also give us the field at points *inside* the sphere. For example, the computation can be used to obtain a good approximation to the field inside an atomic nucleus. In spite of the fact that the protons in a nucleus repel each other, they are, because of the strong nuclear forces, spread nearly uniformly throughout the body of the nucleus.

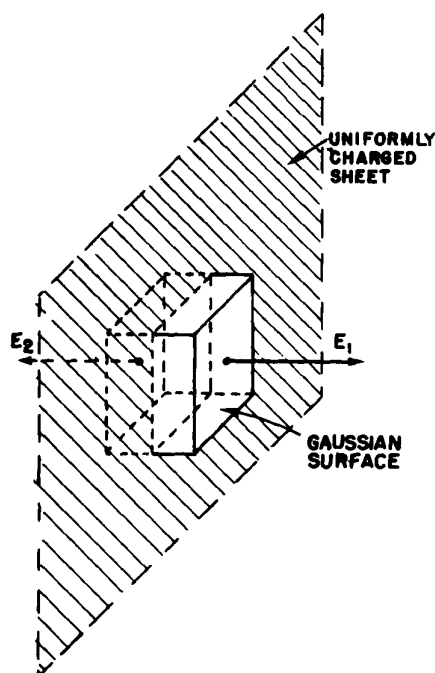


Fig. 5-6. The electric field near a uniformly charged sheet can be found by applying Gauss' law to an imaginary box.

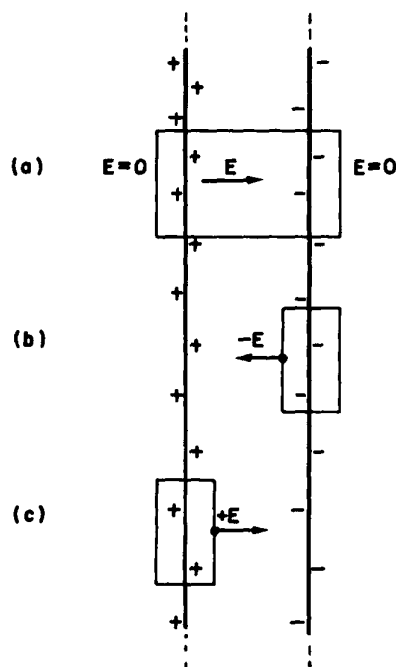


Fig. 5-7. The field between two charged sheets is σ/ϵ_0 .

Suppose that we have a sphere of radius R filled uniformly with charge. Let ρ be the charge per unit volume. Again using arguments of symmetry, we assume the field to be radial and equal in magnitude at all points at the same distance from the center. To find the field at the distance r from the center, we take a spherical gaussian surface of radius r ($r < R$), as shown in Fig. 5-8. The flux out of this surface is

$$4\pi r^2 E.$$

The charge inside our gaussian surface is the volume inside times ρ , or

$$\frac{4}{3}\pi r^3 \rho.$$

Using Gauss' law, it follows that the magnitude of the field is given by

$$E = \frac{\rho r}{3\epsilon_0} \quad (r < R). \quad (5.7)$$

You can see that this formula gives the proper result for $r = R$. The electric field is *proportional* to the radius and is directed radially outward.

The arguments we have just given for a uniformly charged sphere can be applied also to a thin spherical shell of charge. Assuming that the field is everywhere radial and is spherically symmetric, one gets immediately from Gauss' law that the field outside the shell is like that of a point charge, while the field everywhere inside the shell is zero. (A gaussian surface inside the shell will contain no charge.)

5-8 Is the field of a point charge exactly $1/r^2$?

If we look in a little more detail at *how* the field inside the shell gets to be zero, we can see more clearly why it is that Gauss' law is true only because the coulomb force depends exactly on the square of the distance. Consider any point P inside a uniform spherical shell of charge. Imagine a small cone whose apex is at P and which extends to the surface of the sphere, where it cuts out a small surface area Δa_1 , as in Fig. 5-9. An exactly symmetric cone diverging from the opposite side of P would cut out the surface area Δa_2 . If the distances from P to these two elements of area are r_1 and r_2 , the areas are in the ratio

$$\frac{\Delta a_2}{\Delta a_1} = \frac{r_2^2}{r_1^2}.$$

(You can show this by geometry for any point P inside the sphere.)

If the surface of the sphere is uniformly charged, the charge Δq on each of the elements of area is proportional to the area, so

$$\frac{\Delta q_2}{\Delta q_1} = \frac{\Delta a_2}{\Delta a_1}.$$

Coulomb's law then says that the magnitudes of the fields produced at P by these two surface elements are in the ratio

$$\frac{E_2}{E_1} = \frac{q_2/r_2^2}{q_1/r_1^2} = 1.$$

The fields cancel exactly. Since all parts of the surface can be paired off in the same way, the total field at P is zero. But you can see that it would not be so if the exponent of r in Coulomb's law were not exactly two.

The validity of Gauss' law depends upon the inverse square law of Coulomb. If the force law were not exactly the inverse square, it would not be true that the field inside a uniformly charged sphere would be exactly zero. For instance, if the force varied more rapidly, like, say, the inverse cube of r , that portion of the surface which is nearer to an interior point would produce a field which is larger than that which is farther away, resulting in a radial inward field for a positive surface

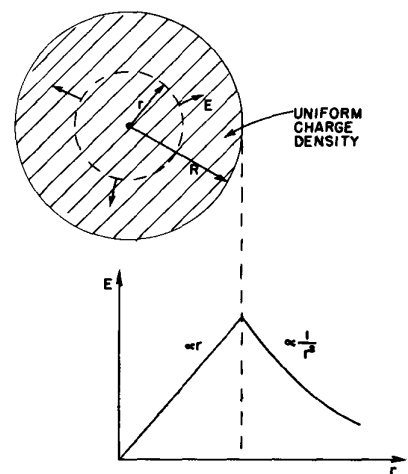


Fig. 5-8. Gauss' law can be used to find the field inside a uniformly charged sphere.

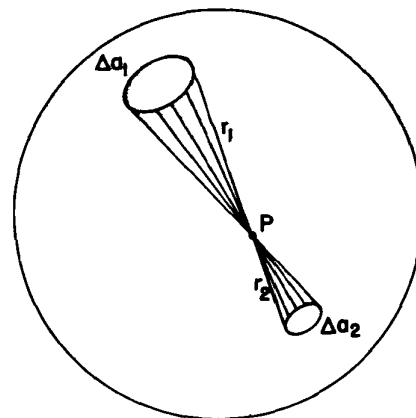


Fig. 5-9. The field is zero at any point P inside a spherical shell of charge.

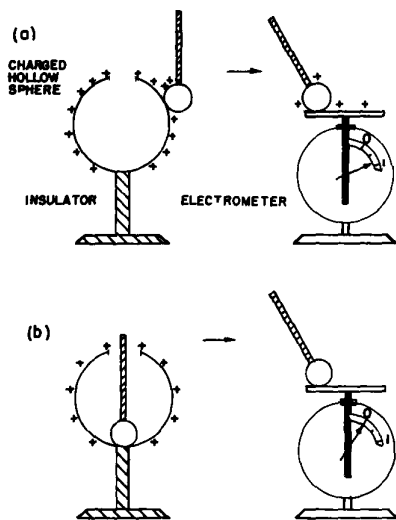


Fig. 5-10. The electric field is zero inside a closed conducting shell.

charge. These conclusions suggest an elegant way of finding out whether the inverse square law is precisely correct. We need only determine whether or not the field inside of a uniformly charged spherical shell is precisely zero.

It is lucky that such a method exists. It is usually difficult to measure a physical quantity to high precision—a one percent result may not be too difficult, but how would one go about measuring, say, Coulomb's law to an accuracy of one part in a billion? It is almost certainly not possible with the best available techniques to measure the *force* between two charged objects with such an accuracy. But by determining only that the electric fields inside a charged sphere are *smaller* than some value we can make a highly accurate measurement of the correctness of Gauss' law, and hence of the inverse square dependence of Coulomb's law. What one does, in effect, is *compare* the force law to an ideal inverse square. Such comparisons of things that are equal, or nearly so, are usually the bases of the most precise physical measurements.

How shall we observe the field inside a charged sphere? One way is to try to charge an object by touching it to the inside of a spherical conductor. You know that if we touch a small metal ball to a charged object and then touch it to an electrometer the meter will become charged and the pointer will move from zero (Fig. 5-10a). The ball picks up charge because there are electric fields outside the charged sphere that cause charges to run onto (or off) the little ball. If you do the same experiment by touching the little ball to the *inside* of the charged sphere, you find that no charge is carried to the electrometer. With such an experiment you can easily show that the field inside is, at most, a few percent of the field outside, and that Gauss' law is at least approximately correct.

It appears that Benjamin Franklin was the first to notice that the field inside a conducting shell is zero. The result seemed strange to him. When he reported his observation to Priestley, the latter suggested that it might be connected with an inverse square law, since it was known that a spherical shell of matter produced no gravitational field inside. But Coulomb didn't measure the inverse square dependence until 18 years later, and Gauss' law came even later still.

Gauss' law has been checked carefully by putting an electrometer inside a large sphere and observing whether any deflections occur when the sphere is charged to a high voltage. A null result is always obtained. Knowing the geometry of the apparatus and the sensitivity of the meter, it is possible to compute the minimum field that would be observed. From this number it is possible to place an upper limit on the deviation of the exponent from two. If we write that the electrostatic force depends on $r^{-2+\epsilon}$, we can place an upper bound on ϵ . By this method Maxwell determined that ϵ was less than 1/10,000. The experiment was repeated and improved upon in 1936 by Plimpton and Laughton. They found that Coulomb's exponent differs from two by less than one part in a billion.

Now that brings up an interesting question: How accurate do we know this Coulomb law to be in various circumstances? The experiments we just described measure the dependence of the field on distance for distances of some tens of centimeters. But what about the distances inside an atom—in the hydrogen atom, for instance, where we believe the electron is attracted to the nucleus by the same inverse square law? It is true that quantum mechanics must be used for the mechanical part of the behavior of the electron, but the force is the usual electrostatic one. In the formulation of the problem, the potential energy of an electron must be known as a function of distance from the nucleus, and Coulomb's law gives a potential which varies inversely with the first power of the distance. How accurately is the exponent known for such small distances? As a result of very careful measurements in 1947 by Lamb and Retherford on the relative positions of the energy levels of hydrogen, we know that the exponent is correct again to one part in a billion on the atomic scale—that is, at distances of the order of one angstrom (10^{-8} centimeter).

The accuracy of the Lamb-Retherford measurement was possible again because of a physical "accident." Two of the states of a hydrogen atom are expected to have almost identical energies *only* if the potential varies exactly as $1/r$. A measurement was made of the very slight *difference* in energies by finding

the frequency ω of the photons that are emitted or absorbed in the transition from one state to the other, using for the energy difference $\Delta E = \hbar\omega$. Computations showed that ΔE would have been noticeably different from what was observed if the exponent in the force law $1/r^2$ differed from 2 by as much as one part in a billion.

Is the same exponent correct at still shorter distances? From measurements in nuclear physics it is found that there are electrostatic forces at typical nuclear distances—at about 10^{-13} centimeter—and that they still vary approximately as the inverse square. We shall look at some of the evidence in a later chapter. Coulomb's law is, we know, still valid, at least to some extent, at distances of the order of 10^{-13} centimeter.

How about 10^{-14} centimeter? This range can be investigated by bombarding protons with very energetic electrons and observing how they are scattered. Results to date seem to indicate that the law fails at these distances. The electrical force seems to be about 10 times too weak at distances less than 10^{-14} centimeter. Now there are two possible explanations. One is that the Coulomb law does not work at such small distances; the other is that our objects, the electrons and protons, are not point charges. Perhaps either the electron or proton, or both, is some kind of a smear. Most physicists prefer to think that the charge of the proton is smeared. We know that protons interact strongly with mesons. This implies that a proton will, from time to time, exist as a neutron with a π^+ meson around it. Such a configuration would act—on the average—like a little sphere of positive charge. We know that the field from a sphere of charge does not vary as $1/r^2$ all the way into the center. It is quite likely that the proton charge is smeared, but the theory of pions is still quite incomplete, so it may also be that Coulomb's law fails at very small distances. The question is still open.

One more point: The inverse square law is valid at distances like one meter and also at 10^{-10} m; but is the coefficient $1/4\pi\epsilon_0$ the same? The answer is yes; at least to an accuracy of 15 parts in a million.

We go back now to an important matter that we slighted when we spoke of the experimental verification of Gauss' law. You may have wondered how the experiment of Maxwell or of Plimpton and Laughton could give such an accuracy unless the spherical conductor they used was a perfect sphere. An accuracy of one part in a billion is really something to achieve, and you might well ask whether they could make a sphere which was that precise. There are certain to be slight irregularities in any real sphere and if there are irregularities, will they not produce fields inside? We wish to show now that it is not necessary to have a perfect sphere. It is possible, in fact, to show that there is no field inside a closed conducting shell of *any* shape. In other words, the experiments depended on $1/r^2$, but had nothing to do with the surface being a sphere (except that with a sphere it is easier to calculate what the fields *would* be if Coulomb had been wrong), so we take up that subject now. To show this, it is necessary to know some of the properties of electrical conductors.

5-9 The fields of a conductor

An electrical conductor is a solid that contains many "free" electrons. The electrons can move around freely *in* the material, but cannot leave the surface. In a metal there are so many free electrons that any electric field will set large numbers of them into motion. Either the current of electrons so set up must be continually kept moving by external sources of energy, or the motion of the electrons will cease as they discharge the sources producing the initial field. In "electrostatic" situations, we do not consider continuous sources of current (they will be considered later when we study magnetostatics), so the electrons move only until they have arranged themselves to produce zero electric field everywhere inside the conductor. (This usually happens in a small fraction of a second.) If there were any field left, this field would urge still more electrons to move; the only electrostatic solution is that the field is everywhere zero inside.

Now consider the *interior* of a charged conducting object. (By "interior" we mean in the *metal* itself.) Since the metal is a conductor, the interior field must

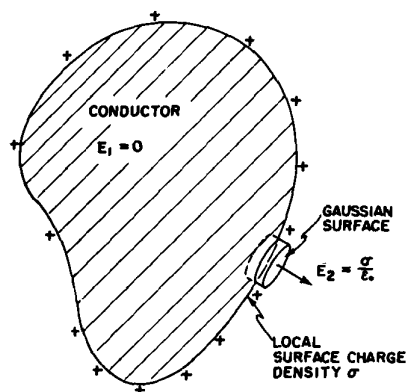


Fig. 5-11. The electric field just outside the surface of a conductor is proportional to the local surface density of charge.

be zero, and so the gradient of the potential ϕ is zero. That means that ϕ does not vary from point to point. Every conductor is an equipotential *region*, and its surface is an equipotential surface. Since in a conducting material the electric field is everywhere zero, the divergence of E is zero, and by Gauss' law the charge density in the *interior* of the conductor must be zero.

If there can be no charges in a conductor, how can it ever be charged? What do we mean when we say a conductor is "charged"? Where are the charges? The answer is that they reside at the surface of the conductor, where there are strong forces to keep them from leaving—they are not completely "free." When we study solid-state physics, we shall find that the excess charge of any conductor is on the average within one or two atomic layers of the surface. For our present purposes, it is accurate enough to say that if any charge is put on, or *in*, a conductor it all accumulates on the surface; there is no charge in the interior of a conductor.

We note also that the electric field *just outside* the surface of a conductor must be normal to the surface. There can be no tangential component. If there were a tangential component, the electrons would move *along* the surface; there are no forces preventing that. Saying it another way: we know that the electric field lines must always go at right angles to an equipotential surface.

We can also, using Gauss' law, relate the field strength just outside a conductor to the local density of the charge at the surface. For a gaussian surface, we take a small cylindrical box half inside and half outside the surface, like the one shown in Fig. 5-11. There is a contribution to the total flux of E only from the side of the box outside the conductor. The field just outside the surface of a conductor is then

Outside a conductor:

$$E = \frac{\sigma}{\epsilon_0}, \quad (5.8)$$

where σ is the *local* surface charge density.

Why does a sheet of charge on a conductor produce a different field than *just* a sheet of charge? In other words, why is (5.8) twice as large as (5.3)? The reason, of course, is that we have *not* said for the conductor that there are no "other" charges around. There must, in fact, be some to make $E = 0$ in the conductor. The charges in the immediate neighborhood of a point P on the surface do, in fact, give a field $E_{\text{local}} = \sigma_{\text{local}}/2\epsilon_0$ both inside and outside the surface. But all the rest of the charges on the conductor "conspire" to produce an additional field at the point P equal in magnitude to E_{local} . The total field inside goes to zero and the field outside to $2E_{\text{local}} = \sigma/\epsilon_0$.

5-10 The field in a cavity of a conductor

We return now to the problem of the hollow container—a conductor with a cavity. There is no field in the *metal*, but what about in the *cavity*? We shall show that if the cavity is *empty* then there are no fields in it, *no matter what the shape* of the conductor or the cavity—say for the one in Fig. 5-12. Consider a gaussian surface, like S in Fig. 5-12, that encloses the cavity but stays everywhere in the conducting material. Everywhere on S the field is zero, so there is no flux through S and the *total* charge inside S is zero. For a spherical shell, one could then argue from symmetry that there could be *no* charge inside. But, in general, we can only say that there are equal amounts of positive and negative charge on the inner surface of the conductor. There *could* be a positive surface charge on one part and a negative one somewhere else, as indicated in Fig. 5-12. Such a thing cannot be ruled out by Gauss' law.

What really happens, of course, is that any equal and opposite charges on the inner surface would slide around to meet each other, cancelling out completely. We can show that they must cancel completely by using the law that the circulation of E is always zero (electrostatics). Suppose there were charges on some parts of the inner surface. We know that there would have to be an equal number of opposite charges somewhere else. Now any lines of E would have to start on the

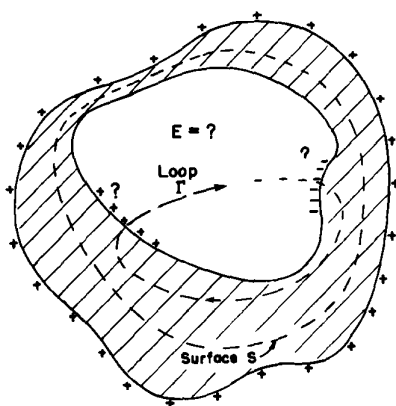


Fig. 5-12. What is the field in an empty cavity of a conductor, for any shape?

positive charges and end on the negative charges (since we are considering only the case that there are no free charges in the cavity). Now imagine a loop Γ that crosses the cavity along a line of force from some positive charge to some negative charge, and returns to its starting point via the conductor (as in Fig. 5-12). The integral along such a line of force from the positive to the negative charges would not be zero. The integral through the metal is zero, since $E = 0$. So we would have

$$\oint E \cdot ds \neq 0???$$

But the line integral of E around any closed loop in an electrostatic field is always zero. So there can be no fields inside the empty cavity, nor any charges on the inside surface.

You should notice carefully one important qualification we have made. We have always said “inside an *empty*” cavity. If some charges are *placed* at some fixed locations in the cavity—as on an insulator or on a small conductor insulated from the main one—then there *can* be fields in the cavity. But then that is not an “empty” cavity.

We have shown that if a cavity is completely enclosed by a conductor, no static distribution of charges *outside* can ever produce any fields inside. This explains the principle of “shielding” electrical equipment by placing it in a metal can. The same arguments can be used to show that no static distribution of charges *inside* a closed conductor can produce any fields *outside*. Shielding works both ways! In electrostatics—but not in varying fields—the fields on the two sides of a closed conducting shell are completely independent.

Now you see why it was possible to check Coulomb’s law to such a great precision. The shape of the hollow shell used doesn’t matter. It doesn’t need to be spherical; it could be square! If Gauss’ law is exact, the field inside is always zero. Now you also understand why it is safe to sit inside the high-voltage terminal of a million-volt van de Graaff generator, without worrying about getting a shock—because of Gauss’ law.

The Electric Field in Various Circumstances

6-1 Equations of the electrostatic potential

This chapter will describe the behavior of the electric field in a number of different circumstances. It will provide some experience with the way the electric field behaves, and will describe some of the mathematical methods which are used to find this field.

We begin by pointing out that the whole mathematical problem is the solution of two equations, the Maxwell equations for electrostatics:

$$\nabla \cdot \mathbf{E} = \frac{\rho}{\epsilon_0}, \quad (6.1)$$

$$\nabla \times \mathbf{E} = 0. \quad (6.2)$$

In fact, the two can be combined into a single equation. From the second equation, we know at once that we can describe the field as the gradient of a scalar (see Section 3-7):

$$\mathbf{E} = -\nabla\phi. \quad (6.3)$$

We may, if we wish, completely describe any particular electric field in terms of its potential ϕ . We obtain the differential equation that ϕ must obey by substituting Eq. (6.3) into (6.1), to get

$$\nabla \cdot \nabla\phi = -\frac{\rho}{\epsilon_0}. \quad (6.4)$$

The divergence of the gradient of ϕ is the same as ∇^2 operating on ϕ :

$$\nabla \cdot \nabla\phi = \nabla^2\phi = \frac{\partial^2\phi}{\partial x^2} + \frac{\partial^2\phi}{\partial y^2} + \frac{\partial^2\phi}{\partial z^2}, \quad (6.5)$$

so we write Eq. (6.4) as

$$\nabla^2\phi = -\frac{\rho}{\epsilon_0}. \quad (6.6)$$

The operator ∇^2 is called the Laplacian, and Eq. (6.6) is called the Poisson equation. The entire subject of electrostatics, from a mathematical point of view, is merely a study of the solutions of the single equation (6.6). Once ϕ is obtained by solving Eq. (6.6) we can find \mathbf{E} immediately from Eq. (6.3).

We take up first the special class of problems in which ρ is given as a function of x, y, z . In that case the problem is almost trivial, for we already know the solution of Eq. (6.6) for the general case. We have shown that if ρ is known at every point, the potential at point (1) is

$$\phi(1) = \int \frac{\rho(2) dV_2}{4\pi\epsilon_0 r_{12}}, \quad (6.7)$$

where $\rho(2)$ is the charge density, dV_2 is the volume element at point (2), and r_{12} is the distance between points (1) and (2). The solution of the *differential* equation (6.6) is reduced to an *integration* over space. The solution (6.7) should be especially noted, because there are many situations in physics that lead to equations like

$$\nabla^2 (\text{something}) = (\text{something else}),$$

and Eq. (6.7) is a prototype of the solution for any of these problems.

The solution of electrostatic field problems is thus completely straightforward when the positions of all the charges are known. Let's see how it works in a few examples.

6-1 Equations of the electrostatic potential

6-2 The electric dipole

6-3 Remarks on vector equations

6-4 The dipole potential as a gradient

6-5 The dipole approximation for an arbitrary distribution

6-6 The fields of charged conductors

6-7 The method of images

6-8 A point charge near a conducting plane

6-9 A point charge near a conducting sphere

6-10 Condensers; parallel plates

6-11 High-voltage breakdown

6-12 The field emission microscope

Review. Chapter 23, Vol. I, Resonance

6-2 The electric dipole

First, take two point charges, $+q$ and $-q$, separated by the distance d . Let the z -axis go through the charges, and pick the origin halfway between, as shown in Fig. 6-1. Then, using (4.24), the potential from the two charges is given by

$$\phi(x, y, z) = \frac{1}{4\pi\epsilon_0} \left[\frac{q}{\sqrt{[z - (d/2)]^2 + x^2 + y^2}} + \frac{-q}{\sqrt{[z + (d/2)]^2 + x^2 + y^2}} \right]. \quad (6.8)$$

We are not going to write out the formula for the electric field, but we can always calculate it once we have the potential. So we have solved the problem of two charges.

There is an important special case in which the two charges are very close together—which is to say that we are interested in the fields only at distances from the charges large in comparison with their separation. We call such a close pair of charges a *dipole*. Dipoles are very common.

A “dipole” antenna can often be approximated by two charges separated by a small distance—if we don’t ask about the field too close to the antenna. (We are usually interested in antennas with *moving* charges; then the equations of statics do not really apply, but for some purposes they are an adequate approximation.)

More important perhaps, are atomic dipoles. If there is an electric field in any material, the electrons and protons feel opposite forces and are displaced relative to each other. In a conductor, you remember, some of the electrons move to the surfaces, so that the field inside becomes zero. In an insulator the electrons cannot move very far; they are pulled back by the attraction of the nucleus. They do, however, shift a little bit. So although an atom, or molecule, remains neutral in an external electric field, there is a very tiny separation of its positive and negative charges and it becomes a microscopic dipole. If we are interested in the fields of these atomic dipoles in the neighborhood of ordinary-sized objects, we are normally dealing with distances large compared with the separations of the pairs of charges.

In some molecules the charges are somewhat separated even in the absence of external fields, because of the form of the molecule. In a water molecule, for example, there is a net negative charge on the oxygen atom and a net positive charge on each of the two hydrogen atoms, which are not placed symmetrically but as in Fig. 6-2. Although the charge of the whole molecule is zero, there is a charge distribution with a little more negative charge on one side and a little more positive charge on the other. This arrangement is certainly not as simple as two point charges, but when seen from far away the system acts like a dipole. As we shall see a little later, the field at large distances is not sensitive to the fine details.

Let’s look, then, at the field of two opposite charges with a small separation d . If d becomes zero, the two charges are on top of each other, the two potentials cancel, and there is no field. But if they are not exactly on top of each other, we can get a good approximation to the potential by expanding the terms of (6.8) in a power series in the small quantity d (using the binomial expansion). Keeping terms only to first order in d , we can write

$$\left(z - \frac{d}{2}\right)^2 \approx z^2 - zd.$$

It is convenient to write

$$x^2 + y^2 + z^2 = r^2.$$

Then

$$\left(z - \frac{d}{2}\right)^2 + x^2 + y^2 \approx r^2 - zd = r^2 \left(1 - \frac{zd}{r^2}\right),$$

and

$$\frac{1}{\sqrt{[z - (d/2)]^2 + x^2 + y^2}} \approx \frac{1}{\sqrt{r^2[1 - (zd/r^2)]}} \approx \frac{1}{r} \left(1 - \frac{zd}{r^2}\right)^{-1/2}.$$

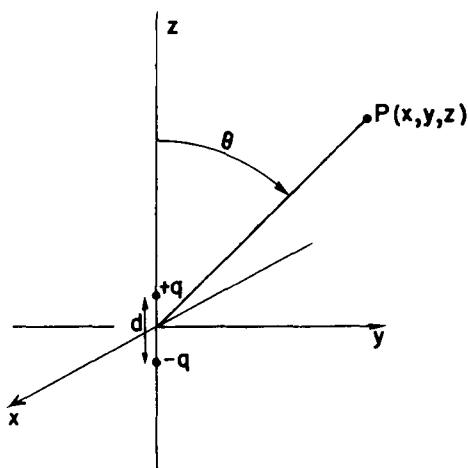


Fig. 6-1. A dipole: two charges $+q$ and $-q$ the distance d apart.

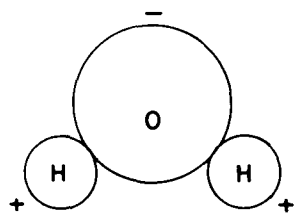


Fig. 6-2. The water molecule H_2O . The hydrogen atoms have slightly less than their share of the electron cloud; the oxygen, slightly more.

Using the binomial expansion again for $[1 - (zd/r^2)]^{-1/2}$ —and throwing away terms with higher powers than the square of d —we get

$$\frac{1}{r} \left(1 + \frac{1}{2} \frac{zd}{r^2} \right).$$

Similarly,

$$\frac{1}{\sqrt{[z + (d/2)]^2 + x^2 + y^2}} \approx \frac{1}{r} \left(1 - \frac{1}{2} \frac{zd}{r^2} \right).$$

The difference of these two terms gives for the potential

$$\phi(x, y, z) = \frac{1}{4\pi\epsilon_0} \frac{z}{r^3} qd. \quad (6.9)$$

The potential, and hence the field, which is its derivative, is proportional to qd , the product of the charge and the separation. This product is defined as the *dipole moment* of the two charges, for which we will use the symbol p (do *not* confuse with momentum!):

$$p = qd. \quad (6.10)$$

Equation (6.9) can also be written as

$$\phi(x, y, z) = \frac{1}{4\pi\epsilon_0} \frac{p \cos \theta}{r^2}, \quad (6.11)$$

since $z/r = \cos \theta$, where θ is the angle between the axis of the dipole and the radius vector to the point (x, y, z) —see Fig. 6-1. The *potential* of a dipole decreases as $1/r^2$ for a given direction from the axis (whereas for a point charge it goes as $1/r$). The electric field E of the dipole will then decrease as $1/r^3$.

We can put our formula into a vector form if we define p as a vector whose magnitude is p and whose direction is along the axis of the dipole, pointing from q_- toward q_+ . Then

$$\cos \theta = p \cdot e_r, \quad (6.12)$$

where e_r is the unit radial vector (Fig. 6-3). We can also represent the point (x, y, z) by r . Then

$$\text{Dipole potential:} \quad \phi(r) = \frac{1}{4\pi\epsilon_0} \frac{p \cdot e_r}{r^2} = \frac{1}{4\pi\epsilon_0} \frac{p \cdot r}{r^3}. \quad (6.13)$$

This formula is valid for a dipole with any orientation and position if r represents the vector from the dipole to the point of interest.

If we want the electric field of the dipole we can get it by taking the gradient of ϕ . For example, the z -component of the field is $-\partial\phi/\partial z$. For a dipole oriented along the z -axis we can use (6.9):

$$-\frac{\partial\phi}{\partial z} = -\frac{p}{4\pi\epsilon_0} \frac{\partial}{\partial z} \left(\frac{z}{r^3} \right) = -\frac{p}{4\pi\epsilon_0} \left(\frac{1}{r^3} - \frac{3z^2}{r^5} \right),$$

or

$$E_z = \frac{p}{4\pi\epsilon_0} \frac{3 \cos^2 \theta - 1}{r^3}. \quad (6.14)$$

The x - and y -components are

$$E_x = \frac{p}{4\pi\epsilon_0} \frac{3zx}{r^5}, \quad E_y = \frac{p}{4\pi\epsilon_0} \frac{3zy}{r^5}.$$

These two can be combined to give one component directed *perpendicular* to the z -axis, which we will call the transverse component E_\perp :

$$E_\perp = \sqrt{E_x^2 + E_y^2} = \frac{p}{4\pi\epsilon_0} \frac{3z}{r^5} \sqrt{x^2 + y^2}$$

or

$$E_\perp = \frac{p}{4\pi\epsilon_0} \frac{3 \cos \theta \sin \theta}{r^3}. \quad (6.15)$$

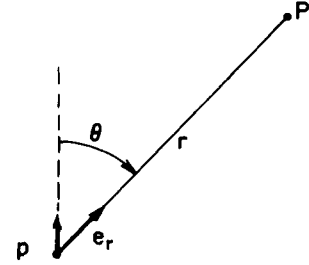


Fig. 6-3. Vector notation for a dipole.

The transverse component E_{\perp} is in the x - y plane and points directly away from the *axis* of the dipole. The total field, of course, is

$$E = \sqrt{E_z^2 + E_{\perp}^2}.$$

The dipole field varies inversely as the cube of the distance from the dipole. On the axis, at $\theta = 0$, it is twice as strong as at $\theta = 90^\circ$. At both of these special angles the electric field has only a z -component, but of opposite sign at the two places (Fig. 6-4).

6-3 Remarks on vector equations

This is a good place to make a general remark about vector analysis. The fundamental proofs can be expressed by elegant equations in a general form, but in making various calculations and analyses it is always a good idea to choose the axes in some convenient way. Notice that when we were finding the potential of a dipole we chose the z -axis along the direction of the dipole, rather than at some arbitrary angle. This made the work much easier. But then we wrote the equations in vector form so that they would no longer depend on any particular coordinate system. After that, we are allowed to choose any coordinate system we wish, knowing that the relation is, in general, true. It clearly doesn't make any sense to bother with an arbitrary coordinate system at some complicated angle when you can choose a neat system for the particular problem—provided that the result can finally be expressed as a vector equation. So by all means take advantage of the fact that vector equations are independent of any coordinate system.

On the other hand, if you are trying to calculate the divergence of a vector, instead of just looking at $\nabla \cdot \mathbf{E}$ and wondering what it is, don't forget that it can always be spread out as

$$\frac{\partial E_x}{\partial x} + \frac{\partial E_y}{\partial y} + \frac{\partial E_z}{\partial z}.$$

If you can then work out the x -, y -, and z -components of the electric field and differentiate them, you will have the divergence. There often seems to be a feeling that there is something inelegant—some kind of defeat involved—in writing out the components; that somehow there ought always to be a way to do everything with the vector operators. There is often no advantage to it. The first time we encounter a particular kind of problem, it usually helps to write out the components to be sure we understand what is going on. There is nothing inelegant about putting numbers into equations, and nothing inelegant about substituting the derivatives for the fancy symbols. In fact, there is often a certain cleverness in doing just that. Of course when you publish a paper in a professional journal it will look better—and be more easily understood—if you can write everything in vector form. Besides, it saves print.

6-4 The dipole potential as a gradient

We would like to point out a rather amusing thing about the dipole formula, Eq. (6.13). The potential can also be written as

$$\phi = -\frac{1}{4\pi\epsilon_0} \mathbf{p} \cdot \nabla \left(\frac{1}{r} \right). \quad (6.16)$$

If you calculate the gradient of $1/r$, you get

$$\nabla \left(\frac{1}{r} \right) = -\frac{\mathbf{r}}{r^3} = -\frac{\mathbf{e}_r}{r^2},$$

and Eq. (6.16) is the same as Eq. (6.13).

How did we think of that? We just remembered that \mathbf{e}_r/r^2 appeared in the formula for the *field* of a point charge, and that the field was the gradient of a *potential* which has a $1/r$ dependence.

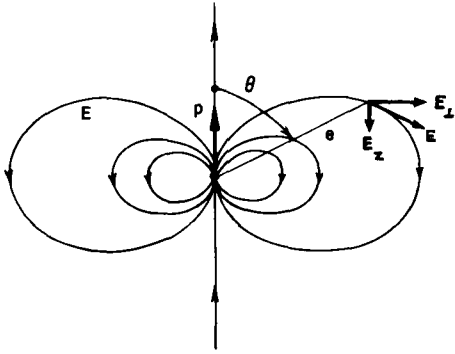


Fig. 6-4. The electric field of a dipole.

There is a *physical* reason for being able to write the dipole potential in the form of Eq. (6.16). Suppose we have a point charge q at the origin. The potential at the point P at (x, y, z) is

$$\phi_0 = \frac{q}{r}.$$

(Let's leave off the $1/4\pi\epsilon_0$ while we make these arguments; we can stick it in at the end.) Now if we move the charge $+q$ up a distance Δz , the potential at P will change a little, by, say, $\Delta\phi_+$. How much is $\Delta\phi_+$? Well, it is just the amount that the potential *would* change if we were to *leave* the charge at the origin and move P *downward* by the same distance Δz (Fig. 6-5). That is,

$$\Delta\phi_+ = -\frac{\partial\phi_0}{\partial z} \Delta z,$$

where by Δz we mean the same as $d/2$. So, using $\phi = q/r$, we have that the potential from the positive charge is

$$\phi_+ = \frac{q}{r} - \frac{\partial}{\partial z} \left(\frac{q}{r} \right) \frac{d}{2}. \quad (6.17)$$

Applying the same reasoning for the potential from the negative charge, we can write

$$\phi_- = \frac{-q}{r} + \frac{\partial}{\partial z} \left(\frac{-q}{r} \right) \frac{d}{2}. \quad (6.18)$$

The total potential is the sum of (6.17) and (6.18):

$$\begin{aligned} \phi = \phi_+ + \phi_- &= -\frac{\partial}{\partial z} \left(\frac{q}{r} \right) d \\ &= -\frac{\partial}{\partial z} \left(\frac{1}{r} \right) qd. \end{aligned} \quad (6.19)$$

For other orientation of the dipole, we could represent the displacement of the positive charge by the vector $\Delta\mathbf{r}_+$. We should then write Eq. (6.17) as

$$\Delta\phi_+ = -\nabla\phi_0 \cdot \Delta\mathbf{r}_+,$$

where $\Delta\mathbf{r}$ is then to be replaced by $d/2$. Completing the derivation as before, Eq. (6.19) would then become

$$\phi = -\nabla \left(\frac{1}{r} \right) \cdot qd.$$

This is the same as Eq. (6.16), if we replace $qd = \mathbf{p}$, and put back the $1/4\pi\epsilon_0$. Looking at it another way, we see that the dipole potential, Eq. (6.13), can be interpreted as

$$\phi = -\mathbf{p} \cdot \nabla\Phi_0, \quad (6.20)$$

where $\Phi_0 = 1/4\pi\epsilon_0 r$ is the potential of a *unit* point charge.

Although we can always find the potential of a known charge distribution by an integration, it is sometimes possible to save time by getting the answer with a clever trick. For example, one can often make use of the superposition principle. If we are given a charge distribution that can be made up of the sum of two distributions for which the potentials are already known, it is easy to find the desired potential by just adding the two known ones. One example of this is our derivation of (6.20), another is the following.

Suppose we have a spherical surface with a distribution of surface charge that varies as the cosine of the polar angle. The integration for this distribution is fairly messy. But, surprisingly, such a distribution can be analyzed by superposition. For imagine a sphere with a uniform *volume* density of positive charge, and another sphere with an equal uniform volume density of negative charge,

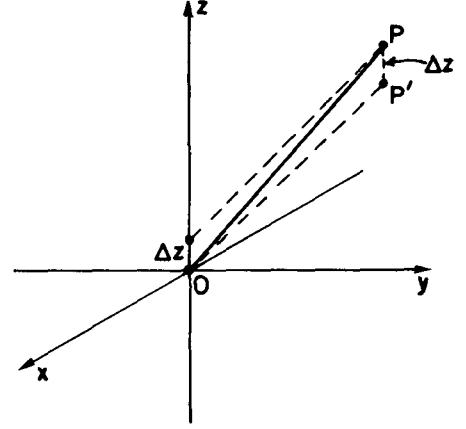
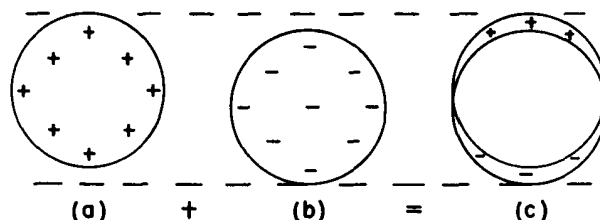


Fig. 6-5. The potential at P from a point charge at Δz above the origin is the same as the potential at P' (Δz below P) from the same charge at the origin.

Fig. 6-6. Two uniformly charged spheres, superposed with a slight displacement, are equivalent to a nonuniform distribution of surface charge.



originally superposed to make a neutral—that is, uncharged—sphere. If the positive sphere is then displaced slightly with respect to the negative sphere, the body of the uncharged sphere would remain neutral, but a little positive charge will appear on one side, and some negative charge will appear on the opposite side, as illustrated in Fig. 6-6. If the relative displacement of the two spheres is small, the net charge is equivalent to a surface charge (on a spherical surface), and the surface charge density will be proportional to the cosine of the polar angle.

Now if we want the potential from this distribution, we do not need to do an integral. We know that the potential from each of the spheres of charge is—for points outside the sphere—the same as from a point charge. The two displaced spheres are like two point charges; the potential is just that of a dipole.

In this way you can show that a charge distribution on a sphere of radius a with a surface charge density

$$\sigma = \sigma_0 \cos \theta$$

produces a field outside the sphere which is just that of a dipole whose moment is

$$p = \frac{4\pi\sigma_0 a^3}{3}.$$

It can also be shown that inside the sphere the field is constant, with the value

$$E = \frac{\sigma_0}{3\epsilon_0}.$$

If θ is the angle from the positive z -axis, the electric field inside the sphere is in the *negative* z -direction. The example we have just considered is not as artificial as it may appear; we will encounter it again in the theory of dielectrics.

6-5 The dipole approximation for an arbitrary distribution

The dipole field appears in another circumstance both interesting and important. Suppose that we have an object that has a complicated distribution of charge—like the water molecule (Fig. 6-2)—and we are interested only in the fields far away. We will show that it is possible to find a relatively simple expression for the fields which is appropriate for distances large compared with the size of the object.

We can think of our object as an assembly of point charges q_i in a certain limited region, as shown in Fig. 6-7. (We can, later, replace q_i by ρdV if we wish.) Let each charge q_i be located at the displacement \mathbf{d}_i from an origin chosen somewhere

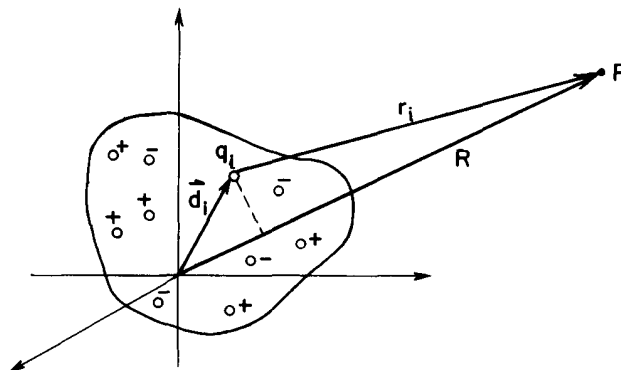


Fig. 6-7. Computation of the potential at a point P at a large distance from a set of charges.

in the middle of the group of charges. What is the potential at the point P , located at \mathbf{R} , where \mathbf{R} is much larger than the maximum \mathbf{d}_i ? The potential from the whole collection is given by

$$\phi = \frac{1}{4\pi\epsilon_0} \sum_i \frac{q_i}{r_i}, \quad (6.21)$$

where r_i is the distance from P to the charge q_i (the length of the vector $\mathbf{R} - \mathbf{d}_i$). Now if the distance from the charges to P , the point of observation, is enormous, each of the r_i 's can be approximated by R . Each term becomes q_i/R , and we can take $1/R$ out as a factor in front of the summation. This gives us the simple result

$$\phi = \frac{1}{4\pi\epsilon_0} \frac{1}{R} \sum q_i = \frac{Q}{4\pi\epsilon_0 R}, \quad (6.22)$$

where Q is just the total charge of the whole object. Thus we find that for points far enough from any lump of charge, the lump looks like a point charge. The result is not too surprising.

But what if there are equal numbers of positive and negative charges? Then the total charge Q of the object is zero. This is not an unusual case; in fact, as we know, objects are usually neutral. The water molecule is neutral, but the charges are not all at one point, so if we are close enough we should be able to see some effects of the separate charges. We need a better approximation than (6.22) for the potential from an arbitrary distribution of charge in a neutral object. Equation (6.21) is still precise, but we can no longer just set $r_i = R$. We need a more accurate expression for r_i . If the point P is at a large distance, r_i will differ from R to an excellent approximation by the projection of \mathbf{d} on \mathbf{R} , as can be seen from Fig. 6-7. (You should imagine that P is really farther away than is shown in the figure.) In other words, if \mathbf{e}_r is the unit vector in the direction of \mathbf{R} , then our next approximation to r_i is

$$r_i \approx R - \mathbf{d}_i \cdot \mathbf{e}_r. \quad (6.23)$$

What we really want is $1/r_i$, which, since $d_i \ll R$, can be written to our approximation as

$$\frac{1}{r_i} \approx \frac{1}{R} \left(1 + \frac{\mathbf{d}_i \cdot \mathbf{e}_r}{R} \right). \quad (6.24)$$

Substituting this in (6.21), we get that the potential is

$$\phi = \frac{1}{4\pi\epsilon_0} \left(\frac{Q}{R} + \sum_i q_i \frac{\mathbf{d}_i \cdot \mathbf{e}_r}{R^2} + \dots \right). \quad (6.25)$$

The three dots indicate the terms of higher order in d/R that we have neglected. These, as well as the ones we have already obtained, are successive terms in a Taylor expansion of $1/r_i$ about $1/R$ in powers of d_i/R .

The first term in (6.25) is what we got before; it drops out if the object is neutral. The second term depends on $1/R^2$, just as for a dipole. In fact, if we *define*

$$\mathbf{p} = \sum q_i \mathbf{d}_i \quad (6.26)$$

as a property of the charge distribution, the second term of the potential (6.25) is

$$\phi = \frac{1}{4\pi\epsilon_0} \frac{\mathbf{p} \cdot \mathbf{e}_r}{R^2}, \quad (6.27)$$

precisely a dipole potential. The quantity \mathbf{p} is called the dipole moment of the distribution. It is a generalization of our earlier definition, and reduces to it for the special case of two point charges.

Our result is that, far enough away from *any* mess of charges that is as a whole neutral, the potential is a dipole potential. It decreases as $1/R^2$ and varies as $\cos \theta$ —and its strength depends on the dipole moment of the distribution of charge. It is for these reasons that dipole fields are important, since the simple case of a pair of point charges is quite rare.

The water molecule, for example, has a rather strong dipole moment. The electric fields that result from this moment are responsible for some of the important properties of water. For many molecules, for example CO_2 , the dipole moment vanishes because of the symmetry of the molecule. For them we should expand still more accurately, obtaining another term in the potential which decreases as $1/R^3$, and which is called a quadrupole potential. We will discuss such cases later.

6-6 The fields of charged conductors

We have now finished with the examples we wish to cover of situations in which the charge distributions is known from the start. It has been a problem without serious complications, involving at most some integrations. We turn now to an entirely new kind of problem, the determination of the fields near charged conductors.

Suppose that we have a situation in which a total charge Q is placed on an arbitrary conductor. Now we will not be able to say exactly where the charges are. They will spread out in some way on the surface. How can we know how the charges have distributed themselves on the surface? They must distribute themselves so that the potential of the surface is constant. If the surface were not an equipotential, there would be an electric field inside the conductor, and the charges would keep moving until it became zero. The general problem of this kind can be solved in the following way. We guess at a distribution of charge and calculate the potential. If the potential turns out to be constant everywhere on the surface, the problem is finished. If the surface is not an equipotential, we have guessed the wrong distribution of charges, and should guess again—hopefully with an improved guess! This can go on forever, unless we are judicious about the successive guesses.

The question of how to guess at the distribution is mathematically difficult. Nature, of course, has time to do it; the charges push and pull until they all balance themselves. When we try to solve the problem, however, it takes us so long to make each trial that that method is very tedious. With an arbitrary group of conductors and charges the problem can be very complicated, and in general it cannot be solved without rather elaborate numerical methods. Such numerical computations, these days, are set up on a computing machine that will do the work for us, once we have told it how to proceed.

On the other hand, there are a lot of little practical cases where it would be nice to be able to find the answer by some more direct method—without having to write a program for a computer. Fortunately, there are a number of cases where the answer can be obtained by squeezing it out of Nature by some trick or other. The first trick we will describe involves making use of solutions we have already obtained for situations in which charges have specified locations.

6-7 The method of images

We have solved, for example, the field of two point charges. Figure 6-8 shows some of the field lines and equipotential surfaces we obtained by the computations in Chapter 5. Now consider the equipotential surface marked *A*. Suppose we were to shape a thin sheet of metal so that it just fits this surface. If we place it right at the surface and adjust its potential to the proper value, no one would ever know it was there, because nothing would be changed.

But notice! We have really solved a *new* problem. We have a situation in which the surface of a curved conductor with a given potential is placed near a point charge. If the metal sheet we placed at the equipotential surface eventually closes on itself (or, in practice, if it goes far enough) we have the kind of situation considered in Section 5-10, in which our space is divided into two regions, one inside and one outside a closed conducting shell. We found there that the fields in the two regions are quite independent of each other. So we would have the same fields outside our curved conductor no matter what is inside. We can even fill up

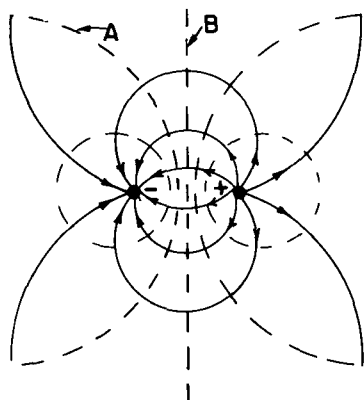


Fig. 6-8. The field lines and equipotentials for two point charges.

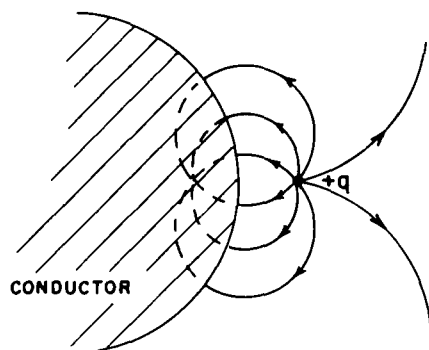


Fig. 6-9. The field outside a conductor shaped like the equipotential *A* of Fig. 6-8.

the whole inside with conducting material. We have found, therefore, the fields for the arrangement of Fig. 6-9. In the space outside the conductor the field is just like that of two point charges, as in Fig. 6-8. Inside the conductor, it is zero. Also—as it must be—the electric field just outside the conductor is normal to the surface.

Thus we can compute the fields in Fig. 6-9 by computing the field due to q and to an imaginary point charge $-q$ at a suitable point. The point charge we “imagine” existing behind the conducting surface is called an *image charge*.

In books you can find long lists of solutions for hyperbolic-shaped conductors and other complicated looking things, and you wonder how anyone ever solved these terrible shapes. They were solved backwards! Someone solved a simple problem with given charges. He then saw that some equipotential surface showed up in a new shape, and he wrote a paper in which he pointed out that the field outside that particular shape can be described in a certain way.

6-8 A point charge near a conducting plane

As the simplest application of the use of this method, let's make use of the plane equipotential surface B of Fig. 6-8. With it, we can solve the problem of a charge in front of a conducting sheet. We just cross out the left-hand half of the picture. The field lines for our solution are shown in Fig. 6-10. Notice that the plane, since it was halfway between the two charges, has zero potential. We have solved the problem of a positive charge next to a grounded conducting sheet.

We have now solved for the total field, but what about the *real* charges that are responsible for it? There are, in addition to our positive point charge, some induced negative charges on the conducting sheet that have been attracted by the positive charge (from large distances away). Now suppose that for some technical reason—or out of curiosity—you would like to know how the negative charges are distributed on the surface. You can find the surface charge density by using the result we worked out in Section 5-6 with Gauss' theorem. The normal com-

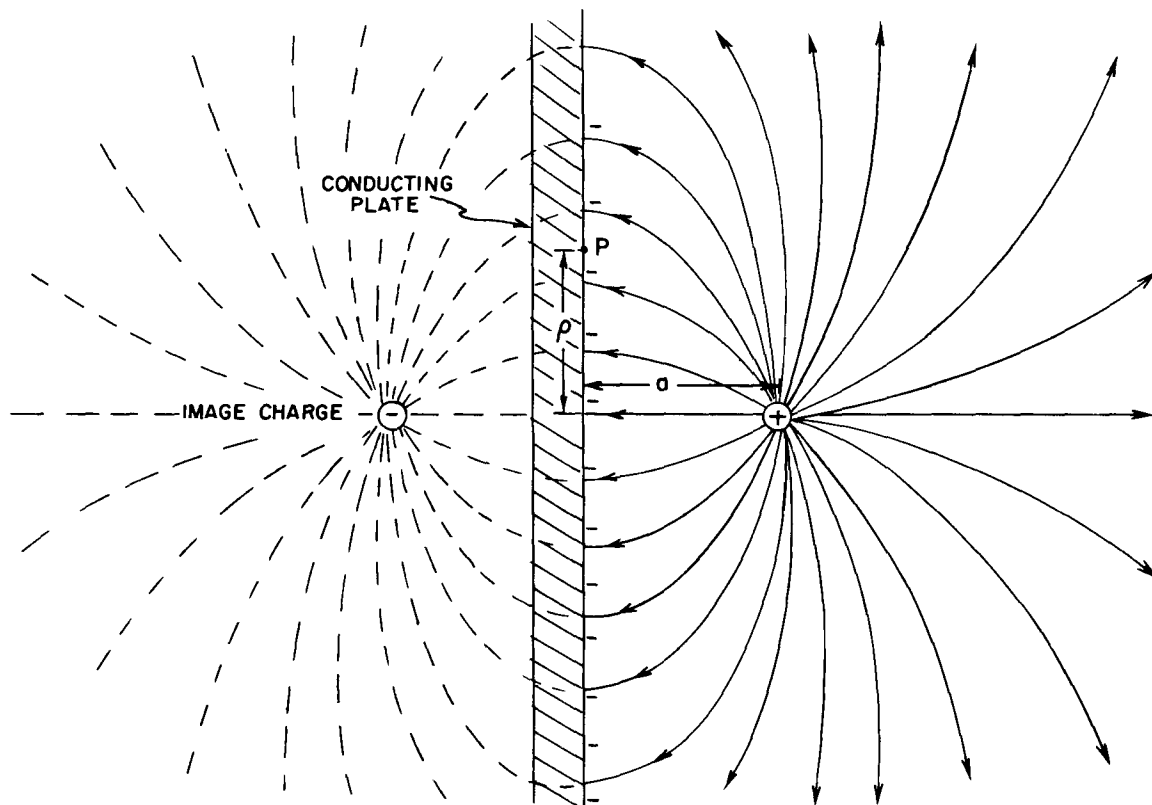


Fig. 6-10. The field of a charge near a plane conducting surface, found by the method of images.

ponent of the electric field just outside a conductor is equal to the density of surface charge σ divided by ϵ_0 . We can obtain the density of charge at any point on the surface by working backwards from the normal component of the electric field at the surface. We know that, because we know the field everywhere.

Consider a point on the surface at the distance ρ from the point directly beneath the positive charge (Fig. 6-10). The electric field at this point is normal to the surface and is directed into it. The component normal to the surface of the field from the *positive* point charge is

$$E_{n+} = -\frac{1}{4\pi\epsilon_0} \frac{aq}{(a^2 + \rho^2)^{3/2}} \quad (6.28)$$

To this we must add the electric field produced by the negative image charge. That just doubles the normal component (and cancels all others), so the charge density σ at any point on the surface is

$$\sigma(\rho) = \epsilon_0 E(\rho) = -\frac{2aq}{4\pi(a^2 + \rho^2)^{3/2}} \quad (6.29)$$

An interesting check on our work is to integrate σ over the whole surface. We find that the total induced charge is $-q$, as it should be.

One further question: Is there a force on the point charge? Yes, because there is an attraction from the induced negative surface charge on the plate. Now that we know what the surface charges are (from Eq. (6.29)), we could compute the force on our positive point charge by an integral. But we also know that the force acting on the positive charge is exactly the same as it *would be* with the negative image charge instead of the plate, because the fields in the neighborhood are the same in both cases. The point charge feels a force toward the plate whose magnitude is

$$F = \frac{1}{4\pi\epsilon_0} \frac{q^2}{(2a)^2} \quad (6.30)$$

We have found the force much more easily than by integrating over all the negative charges.

6-9 A point charge near a conducting sphere

What other surfaces besides a plane have a simple solution? The next most simple shape is a sphere. Let's find the fields around a metal sphere which has a point charge q near it, as shown in Fig. 6-11. Now we must look for a simple physical situation which gives a sphere for an equipotential surface. If we look around at problems people have already solved, we find that someone has noticed that the field of two *unequal* point charges has an equipotential that is a sphere. Aha! If we choose the location of an image charge—and pick the right amount of charge—maybe we can make the equipotential surface fit our sphere. Indeed, it can be done with the following prescription.

Assume that you want the equipotential surface to be a sphere of radius a with its center at the distance b from the charge q . Put an image charge of strength $q' = -q(a/b)$ on the line from the charge to the center of the sphere, and at a distance a^2/b from the center. The sphere will be at zero potential.

The mathematical reason stems from the fact that a sphere is the locus of all points for which the distances from two points are in a constant ratio. Referring to Fig. 6-11, the potential at P from q and q' is proportional to

$$\frac{q}{r_1} + \frac{q'}{r_2}$$

The potential will thus be zero at all points for which

$$\frac{q'}{r_2} = -\frac{q}{r_1} \quad \text{or} \quad \frac{r_2}{r_1} = -\frac{q'}{q}$$

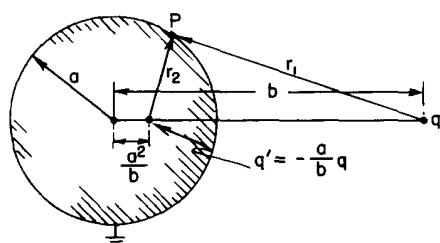


Fig. 6-11. The point charge q induces charges on a grounded conducting sphere whose fields are those of an image charge q' placed at the point shown.

If we place q' at the distance a^2/b from the center, the ratio r_2/r_1 has the constant value a/b . Then if

$$\frac{q'}{q} = -\frac{a}{b}, \quad (6.31)$$

the sphere is an equipotential. Its potential is, in fact, zero.

What happens if we are interested in a sphere that is not at zero potential? That would be so only if its total charge happens accidentally to be q' . Of course if it is grounded, the charges induced on it would have to be just that. But what if it is insulated, and we have put no charge on it? Or if we know that the total charge Q has been put on it? Or just that it has a given potential *not* equal to zero? All these questions are easily answered. We can always add a point charge q'' at the center of the sphere. The sphere still remains an equipotential by superposition; only the magnitude of the potential will be changed.

If we have, for example, a conducting sphere which is initially uncharged and insulated from everything else, and we bring near to it the positive point charge q , the total charge of the sphere will remain zero. The solution is found by using an image charge q' as before, but, in addition, adding a charge q'' at the center of the sphere, choosing

$$q'' = -q' = \frac{a}{b} q. \quad (6.32)$$

The fields everywhere outside the sphere are given by the superposition of the fields of q , q' , and q'' . The problem is solved.

We can see now that there will be a force of attraction between the sphere and the point charge q . It is not zero even though there is no charge on the neutral sphere. Where does the attraction come from? When you bring a positive charge up to a conducting sphere, the positive charge attracts negative charges to the side closer to itself and leaves positive charges on the surface of the far side. The attraction by the negative charges exceeds the repulsion from the positive charges, there is a net attraction. We can find out how large the attraction is by computing the force on q in the field produced by q' and q'' . The total force is the sum of the attractive force between q and a charge $q' = -(a/b)q$, at the distance $b - (a^2/b)$, and the repulsive force between q and a charge $q'' = +(a/b)q$ at the distance b .

Those who were entertained in childhood by the baking powder box which has on its label a picture of a baking powder box which has . . . may be interested in the following problem. Two equal spheres, one with a total charge of $+Q$ and the other with a total charge of $-Q$, are placed at some distance from each other. What is the force between them? The problem can be solved with an infinite number of images. One first approximates each sphere by a charge at its center. These charges will have image charges in the other sphere. The image charges will have images, etc., etc., etc. The solution is like the picture on the box of baking powder—and it converges pretty fast.

6-10 Condensers; parallel plates

We take up now another kind of a problem involving conductors. Consider two large metal plates which are parallel to each other and separated by a distance small compared with their width. Let's suppose that equal and opposite charges have been put on the plates. The charges on each plate will be attracted by the charges on the other plate, and the charges will spread out uniformly on the inner surfaces of the plates. The plates will have surface charge densities $+\sigma$ and $-\sigma$, respectively, as in Fig. 6-12. From Chapter 5 we know that the field between the plates is σ/ϵ_0 , and that the field outside the plates is zero. The plates will have different potentials ϕ_1 and ϕ_2 . For convenience we will call the difference V ; it is often called the "voltage":

$$\phi_1 - \phi_2 = V.$$

(You will find that sometimes people use V for the potential, but we have chosen to use ϕ .)

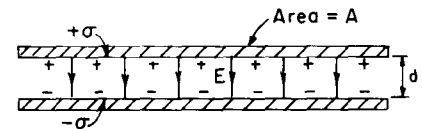


Fig. 6-12 A parallel-plate condenser.

The potential difference V is the work per unit charge required to carry a small charge from one plate to the other, so that

$$V = Ed = \frac{\sigma}{\epsilon_0} d = \frac{d}{\epsilon_0 A} Q, \quad (6.33)$$

where $\pm Q$ is the total charge on each plate, A is the area of the plates, and d is the separation.

We find that the voltage is proportional to the charge. Such a proportionality between V and Q is found for any two conductors in space if there is a plus charge on one and an equal minus charge on the other. The potential difference between them—that is, the voltage—will be proportional to the charge. (We are assuming that there are no other charges around.)

Why this proportionality? Just the superposition principle. Suppose we know the solution for one set of charges, and then we superimpose two such solutions. The charges are doubled, the fields are doubled, and the work done in carrying a unit charge from one point to the other is also doubled. Therefore the potential difference between any two points is proportional to the charges. In particular, the potential difference between the two conductors is proportional to the charges on them. Someone originally wrote the equation of proportionality the other way. That is, they wrote

$$Q = CV,$$

where C is a constant. This coefficient of proportionality is called the *capacity*, and such a system of two conductors is called a *condenser*.* For our parallel-plate condenser

$$C = \frac{\epsilon_0 A}{d} \quad (\text{parallel plates}). \quad (6.34)$$

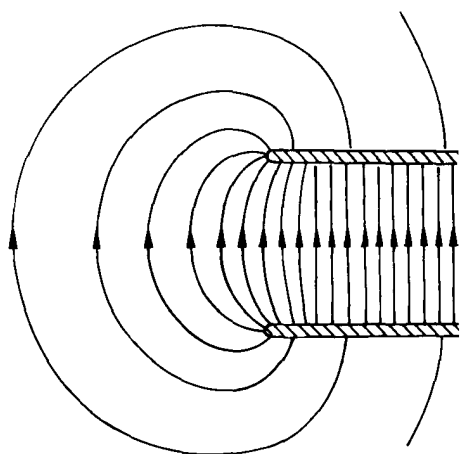


Fig. 6-13. The electric field near the edge of two parallel plates.

This formula is not exact, because the field is not really uniform everywhere between the plates, as we assumed. The field does not just suddenly quit at the edges, but really is more as shown in Fig 6-13. The total charge is not σA , as we have assumed—there is a little correction for the effects at the edges. To find out what the correction is, we will have to calculate the field more exactly and find out just what does happen at the edges. That is a complicated mathematical problem which can, however, be solved by techniques which we will not describe now. The result of such calculations is that the charge density rises somewhat near the edges of the plates. This means that the capacity of the plates is a little higher than we computed. [A very good approximation for the capacity is obtained if we use Eq. (6.34) but take for A the area one *would* get if the plates were extended artificially by a distance $3/8$ of the separation between the plates.]

We have talked about the capacity for two conductors only. Sometimes people talk about the capacity of a single object. They say, for instance, that the capacity of a sphere of radius a is $4\pi\epsilon_0 a$. What they imagine is that the other terminal is another sphere of infinite radius—that when there is a charge $+Q$ on the sphere, the opposite charge, $-Q$, is on an infinite sphere. One can also speak of capacities when there are three or more conductors, a discussion we shall, however, defer.

Suppose that we wish to have a condenser with a very large capacity. We could get a large capacity by taking a very big area and a very small separation. We could put waxed paper between sheets of aluminum foil and roll it up. (If we seal it in plastic, we have a typical radio-type condenser.) What good is it? It is good for storing charge. If we try to store charge on a ball, for example, its potential rises rapidly as we charge it up. It may even get so high that the charge begins to escape into the air by way of sparks. But if we put the same charge on a condenser whose capacity is very large, the voltage developed across the condenser will be small.

* Some people think the words “capacitance” and “capacitor” should be used, instead of “capacity” and “condensor.” We have decided to use the older terminology, because it is still more commonly heard in the physics laboratory—even if not in textbooks!

In many applications in electronic circuits, it is useful to have something which can absorb or deliver large quantities of charge without changing its potential much. A condenser (or “capacitor”) does just that. There are also many applications in electronic instruments and in computers where a condenser is used to get a specified change in voltage in response to a particular change in charge. We have seen a similar application in Chapter 23, Vol. I, where we described the properties of resonant circuits.

From the definition of C , we see that its unit is one coul/volt. This unit is also called a *farad*. Looking at Eq. (6.34), we see that one can express the units of ϵ_0 as farad/meter, which is the unit most commonly used. Typical sizes of condensers run from one micro-microfarad (= 1 picofarad) to millifarads. Small condensers of a few picofarads are used in high-frequency tuned circuits, and capacities up to hundreds or thousands of microfarads are found in power-supply filters. A pair of plates one square centimeter in area with a one millimeter separation have a capacity of roughly one micro-microfarad.

6-11 High-voltage breakdown

We would like now to discuss qualitatively some of the characteristics of the fields around conductors. If we charge a conductor that is not a sphere, but one that has on it a point or a very sharp end, as, for example, the object sketched in Fig. 6-14, the field around the point is much higher than the field in the other regions. The reason is, qualitatively, that charges try to spread out as much as possible on the surface of a conductor, and the tip of a sharp point is as far away as it is possible to be from most of the surface. Some of the charges on the plate get pushed all the way to the tip. A relatively small *amount* of charge on the tip can still provide a large surface *density*; a high charge density means a high field just outside.

One way to see that the field is highest at those places on a conductor where the radius of curvature is smallest is to consider the combination of a big sphere and a little sphere connected by a wire, as shown in Fig. 6-15. It is a somewhat idealized version of the conductor of Fig. 6-14. The wire will have little influence on the fields outside; it is there to keep the spheres at the same potential. Now, which ball has the biggest field at its surface? If the ball on the left has the radius a and carries a charge Q , its potential is about

$$\phi_1 = \frac{1}{4\pi\epsilon_0} \frac{Q}{a}.$$

(Of course the presence of one ball changes the charge distribution on the other, so that the charges are not really spherically symmetric on either. But if we are interested only in an estimate of the fields, we can use the potential of a spherical charge.) If the smaller ball, whose radius is b , carries the charge q , its potential is about

$$\phi_2 = \frac{1}{4\pi\epsilon_0} \frac{q}{b}.$$

But $\phi_1 = \phi_2$, so

$$\frac{Q}{a} = \frac{q}{b}.$$

On the other hand, the field at the surface (see Eq. 5.8) is proportional to the surface charge density, which is like the total charge over the radius squared. We get that

$$\frac{E_a}{E_b} = \frac{Q/a^2}{q/b^2} = \frac{b}{a}. \quad (6.35)$$

Therefore the field is higher at the surface of the small sphere. The fields are in the inverse proportion of the radii.

This result is technically very important, because air will break down if the electric field is too great. What happens is that a loose charge (electron, or ion) somewhere in the air is accelerated by the field, and if the field is very great, the charge can pick up enough speed before it hits another atom to be able to knock an

$$\epsilon_0 \approx \frac{1}{36\pi \times 10^9} \text{ farad/meter}$$

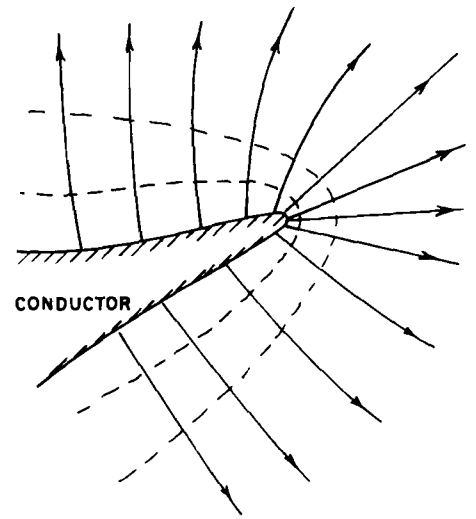


Fig. 6-14. The electric field near a sharp point on a conductor is very high.

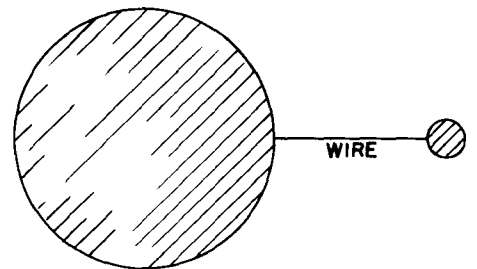


Fig. 6-15. The field of a pointed object can be approximated by that of two spheres at the same potential.

electron off that atom. As a result, more and more ions are produced. Their motion constitutes a discharge, or spark. If you want to charge an object to a high potential and not have it discharge itself by sparks in the air, you must be sure that the surface is smooth, so that there is no place where the field is abnormally large.

6-12 The field-emission microscope

There is an interesting application of the extremely high electric field which surrounds any sharp protuberance on a charged conductor. The *field-emission microscope* depends for its operation on the high fields produced at a sharp metal point.* It is built in the following way. A very fine needle, with a tip whose diameter is about 1000 angstroms, is placed at the center of an evacuated glass sphere (Fig. 6-16.) The inner surface of the sphere is coated with a thin conducting layer of fluorescent material, and a very high potential difference is applied between the fluorescent coating and the needle.

Let's first consider what happens when the needle is negative with respect to the fluorescent coating. The field lines are highly concentrated at the sharp point. The electric field can be as high as 40 million volts per centimeter. In such intense fields, electrons are pulled out of the surface of the needle and accelerated across the potential difference between the needle and the fluorescent layer. When they arrive there they cause light to be emitted, just as in a television picture tube.

The electrons which arrive at a given point on the fluorescent surface are, to an excellent approximation, those which leave the other end of the radial field line, because the electrons will travel along the field line passing from the point to the surface. Thus we see on the surface some kind of an image of the tip of the needle. More precisely, we see a picture of the *emissivity* of the surface of the needle—that is the ease with which electrons can leave the surface of the metal tip. If the resolution were high enough, one could hope to resolve the positions of the individual atoms on the tip of the needle. With electrons, this resolution is not possible for the following reasons. First, there is quantum-mechanical diffraction of the electron waves which blurs the image. Second, due to the internal motions of the electrons in the metal they have a small sideways initial velocity when they leave the needle, and this random transverse component of the velocity causes some smearing of the image. The combination of these two effects limits the resolution to 25 Å or so.

If, however, we reverse the polarity and introduce a small amount of helium gas into the bulb, much higher resolutions are possible. When a helium atom collides with the tip of the needle, the intense field there strips an electron off the helium atom, leaving it positively charged. The helium ion is then accelerated outward along a field line to the fluorescent screen. Since the helium ion is so much heavier than an electron, the quantum-mechanical wavelengths are much smaller. If the temperature is not too high, the effect of the thermal velocities is also smaller than in the electron case. With less smearing of the image a much sharper picture of the point is obtained. It has been possible to obtain magnifications up to 2,000,000 times with the positive ion field-emission microscope—a magnification ten times better than is obtained with the best electron microscope.

Figure 6-17 is an example of the results which were obtained with a field-ion microscope, using a tungsten needle. The center of a tungsten atom ionizes a helium atom at a slightly different rate than the spaces between the tungsten atoms. The pattern of spots on the fluorescent screen shows the arrangement of the *individual atoms* on the tungsten tip. The reason the spots appear in rings can be understood by visualizing a large box of balls packed in a rectangular array, representing the atoms in the metal. If you cut an approximately spherical section out of this box, you will see the ring pattern characteristic of the atomic structure. The field-ion microscope provided human beings with the means of seeing atoms for the first time. This is a remarkable achievement, considering the simplicity of the instrument.

* See E. W. Mueller: "The field-ion microscope," *Advances in Electronics and Electron Physics*, 13, 83-179 (1960). Academic Press, New York

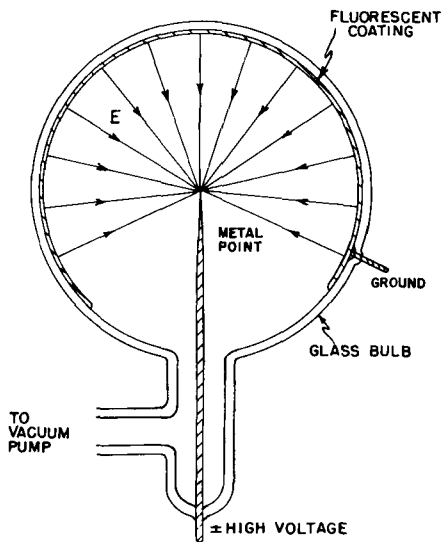


Fig. 6-16. Field-emission microscope.

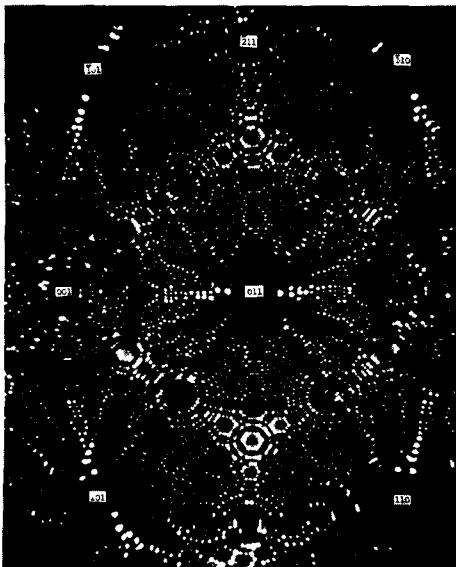


Fig. 6-17. Image produced by a field-emission microscope. [Courtesy of Erwin W. Mueller, Research Prof. of Physics, Pennsylvania State University.]

The Electric Field in Various Circumstances (Continued)

7-1 Methods for finding the electrostatic field

This chapter is a continuation of our consideration of the characteristics of electric fields in various particular situations. We shall first describe some of the more elaborate methods for solving problems with conductors. It is not expected that these more advanced methods can be mastered at this time. Yet it may be of interest to have some idea about the kinds of problems that can be solved, using techniques that may be learned in more advanced courses. Then we take up two examples in which the charge distribution is neither fixed nor is carried by a conductor, but instead is determined by some other law of physics.

As we found in Chapter 6, the problem of the electrostatic field is fundamentally simple when the distribution of charges is specified; it requires only the evaluation of an integral. When there are conductors present, however, complications arise because the charge distribution on the conductors is not initially known; the charge must distribute itself on the surface of the conductor in such a way that the conductor is an equipotential. The solution of such problems is neither direct nor simple.

We have looked at an indirect method of solving such problems, in which we find the equipotentials for some specified charge distribution and replace one of them by a conducting surface. In this way we can build up a catalog of special solutions for conductors in the shapes of spheres, planes, etc. The use of images, described in Chapter 6, is an example of an indirect method. We shall describe another in this chapter.

If the problem to be solved does not belong to the class of problems for which we can construct solutions by the indirect method, we are forced to solve the problem by a more direct method. The mathematical problem of the direct method is the solution of Laplace's equation,

$$\nabla^2\phi = 0, \quad (7.1)$$

subject to the condition that ϕ is a suitable constant on certain boundaries—the surfaces of the conductors. Problems which involve the solution of a differential field equation subject to certain *boundary conditions* are called *boundary-value* problems. They have been the object of considerable mathematical study. In the case of conductors having complicated shapes, there are no general analytical methods. Even such a simple problem as that of a charged cylindrical metal can closed at both ends—a beer can—presents formidable mathematical difficulties. It can be solved only approximately, using numerical methods. The *only* general methods of solution are numerical.

There are a few problems for which Eq. (7.1) can be solved directly. For example, the problem of a charged conductor having the shape of an ellipsoid of revolution can be solved exactly in terms of known special functions. The solution for a thin disc can be obtained by letting the ellipsoid become infinitely oblate. In a similar manner, the solution for a needle can be obtained by letting the ellipsoid become infinitely prolate. However, it must be stressed that the only direct methods of general applicability are the numerical techniques.

Boundary-value problems can also be solved by measurements of a physical analog. Laplace's equation arises in many different physical situations: in steady-state heat flow, in irrotational fluid flow, in current flow in an extended medium,

7-1 Methods for finding the electrostatic field

7-2 Two-dimensional fields; functions of the complex variable

7-3 Plasma oscillations

7-4 Colloidal particles in an electrolyte

7-5 The electrostatic field of a grid

and in the deflection of an elastic membrane. It is frequently possible to set up a physical model which is analogous to an electrical problem which we wish to solve. By the measurement of a suitable analogous quantity on the model, the solution to the problem of interest can be determined. An example of the analog technique is the use of the electrolytic tank for the solution of two-dimensional problems in electrostatics. This works because the differential equation for the potential in a uniform conducting medium is the same as it is for a vacuum.

There are many physical situations in which the variations of the physical fields in one direction are zero, or can be neglected in comparison with the variations in the other two directions. Such problems are called two-dimensional; the field depends on two coordinates only. For example, if we place a long charged wire along the z -axis, then for points not too far from the wire the electric field depends on x and y , but not on z ; the problem is two-dimensional. Since in a two-dimensional problem $\partial/\partial z = 0$, the equation for ϕ in free space is

$$\frac{\partial^2 \phi}{\partial x^2} + \frac{\partial^2 \phi}{\partial y^2} = 0. \quad (7.2)$$

Because the two-dimensional equation is comparatively simple, there is a wide range of conditions under which it can be solved analytically. There is, in fact, a very powerful indirect mathematical technique which depends on a theorem from the mathematics of functions of a complex variable, and which we will now describe.

7-2 Two-dimensional fields; functions of the complex variable

The complex variable z is defined as

$$z = x + iy.$$

(Do not confuse z with the z -coordinate, which we ignore in the following discussion because we assume there is no z -dependence of the fields.) Every point in x and y then corresponds to a complex number z . We can use z as a single (complex) variable, and with it write the usual kinds of mathematical functions $F(z)$. For example,

$$F(z) = z^2,$$

or

$$F(z) = 1/z^3,$$

or

$$F(z) = z \log z,$$

and so forth.

Given any particular $F(z)$ we can substitute $z = x + iy$, and we have a function of x and y —with real and imaginary parts. For example,

$$z^2 = (x + iy)^2 = x^2 - y^2 + 2ixy. \quad (7.3)$$

Any function $F(z)$ can be written as a sum of a pure real part and a pure imaginary part, each part a function of x and y :

$$F(z) = U(x, y) + iV(x, y), \quad (7.4)$$

where $U(x, y)$ and $V(x, y)$ are real functions. Thus from any complex function $F(z)$ two new functions $U(x, y)$ and $V(x, y)$ can be derived. For example, $F(z) = z^2$ gives us the two functions

$$U(x, y) = x^2 - y^2, \quad (7.5)$$

and

$$V(x, y) = 2xy. \quad (7.6)$$

Now we come to a miraculous mathematical theorem which is so delightful that we shall leave a proof of it for one of your courses in mathematics. (We should not reveal all the mysteries of mathematics, or that subject matter would

become too dull.) It is this. For any “ordinary function” (mathematicians will define it better) the functions U and V *automatically* satisfy the relations

$$\frac{\partial U}{\partial x} = \frac{\partial V}{\partial y}, \quad (7.7)$$

$$\frac{\partial V}{\partial x} = -\frac{\partial U}{\partial y}. \quad (7.8)$$

It follows immediately that each of the functions U and V satisfy Laplace’s equation:

$$\frac{\partial^2 U}{\partial x^2} + \frac{\partial^2 U}{\partial y^2} = 0, \quad (7.9)$$

$$\frac{\partial^2 V}{\partial x^2} + \frac{\partial^2 V}{\partial y^2} = 0. \quad (7.10)$$

These equations are clearly true for the functions of (7.5) and (7.6).

Thus, starting with any ordinary function, we can arrive at two functions $U(x, y)$ and $V(x, y)$, which are both solutions of Laplace’s equation in two dimensions. Each function represents a possible electrostatic potential. We can pick *any* function $F(\partial)$ and it should represent *some* electric field problem—in fact, *two* problems, because U and V *each* represent solutions. We can write down as many solutions as we wish—by just making up functions—then we just have to find the *problem* that goes with each solution. It may sound backwards, but it’s a possible approach.

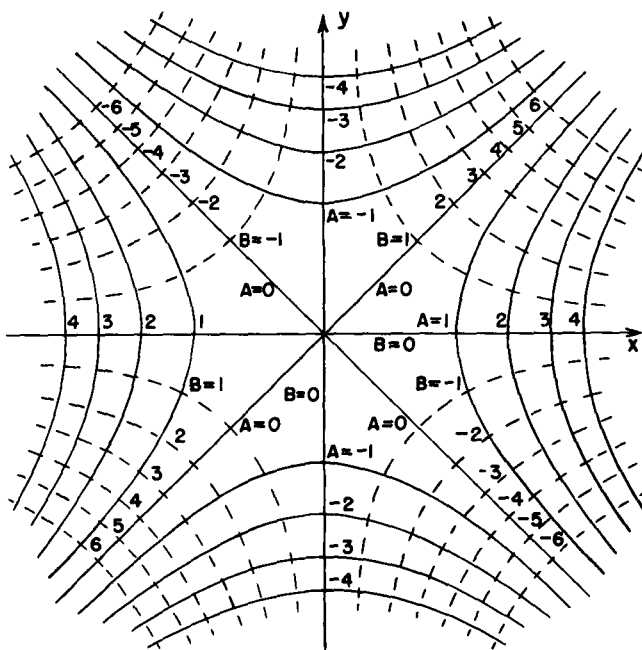


Fig. 7-1. Two sets of orthogonal curves which can represent equipotentials in a two-dimensional electrostatic field.

As an example, let’s see what physics the function $F(\partial) = \partial^2$ gives us. From it we get the two potential functions of (7.5) and (7.6). To see what problem the function U belongs to, we solve for the equipotential surfaces by setting $U = A$, a constant:

$$x^2 - y^2 = A.$$

This is the equation of a rectangular hyperbola. For various values of A , we get the hyperbolas shown in Fig. 7-1. When $A = 0$, we get the special case of diagonal straight lines through the origin.

Such a set of equipotentials corresponds to several possible physical situations. First, it represents the fine details of the field near the point halfway between two

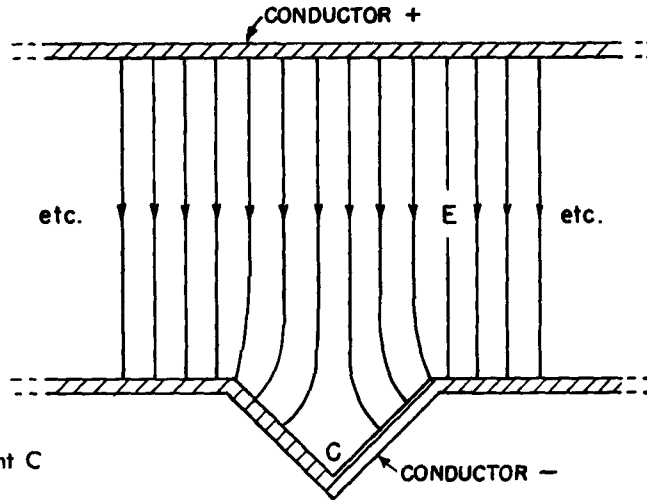


Fig. 7-2. The field near the point C is the same as that in Fig. 7-1.

equal point charges. Second, it represents the field at an inside right-angle corner of a conductor. If we have two electrodes shaped like those in Fig. 7-2, which are held at different potentials, the field near the corner marked C will look just like the field above the origin in Fig. 7-1. The solid lines are the equipotentials, and the broken lines at right angles correspond to lines of E . Whereas at points or protuberances the electric field tends to be high, it tends to be *low* in dents or hollows.

The solution we have found also corresponds to that for a hyperbola-shaped electrode near a right-angle corner, or for two hyperbolas at suitable potentials. You will notice that the field of Fig. 7-1 has an interesting property. The x -component of the electric field, E_x , is given by

$$E_x = -\frac{\partial\phi}{\partial x} = -2x.$$

The electric field is proportional to the distance from the axis. This fact is used to make devices (called quadrupole lenses) that are useful for focusing particle beams (see Section 29-9). The desired field is usually obtained by using four hyperbola-shaped electrodes, as shown in Fig. 7-3. For the electric field lines in Fig. 7-3, we have simply copied from Fig. 7-1 the set of broken-line curves that represent $V = \text{constant}$. We have a bonus! The curves for $V = \text{constant}$ are orthogonal to the ones for $U = \text{constant}$ because of the equations (7.7) and (7.8). Whenever we choose a function $F(\bar{z})$, we get from U and V both the equipotentials and field lines. And you will remember that we have solved either of two problems, depending on which set of curves we call the equipotentials.

As a second example, consider the function

$$F(\bar{z}) = \sqrt{\bar{z}}. \quad (7.11)$$

If we write

$$\bar{z} = x + iy = \rho e^{i\theta},$$

where

$$\rho = \sqrt{x^2 + y^2}$$

and

$$\tan \theta = y/x,$$

then

$$\begin{aligned} F(\bar{z}) &= \rho^{1/2} e^{i\theta/2} \\ &= \rho^{1/2} \left(\cos \frac{\theta}{2} + i \sin \frac{\theta}{2} \right), \end{aligned}$$

from which

$$F(\bar{z}) = \left[\frac{(x^2 + y^2)^{1/2} + x}{2} \right]^{1/2} + i \left[\frac{(x^2 + y^2)^{1/2} - x}{2} \right]^{1/2}. \quad (7.12)$$

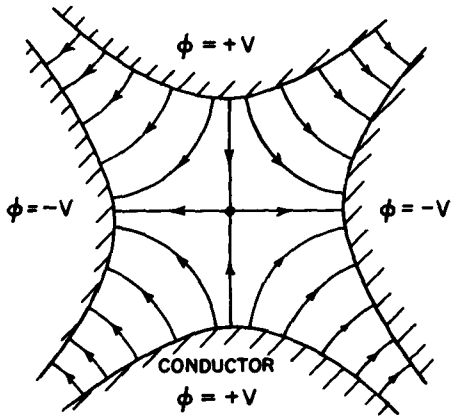


Fig. 7-3. The field in a quadrupole lens.

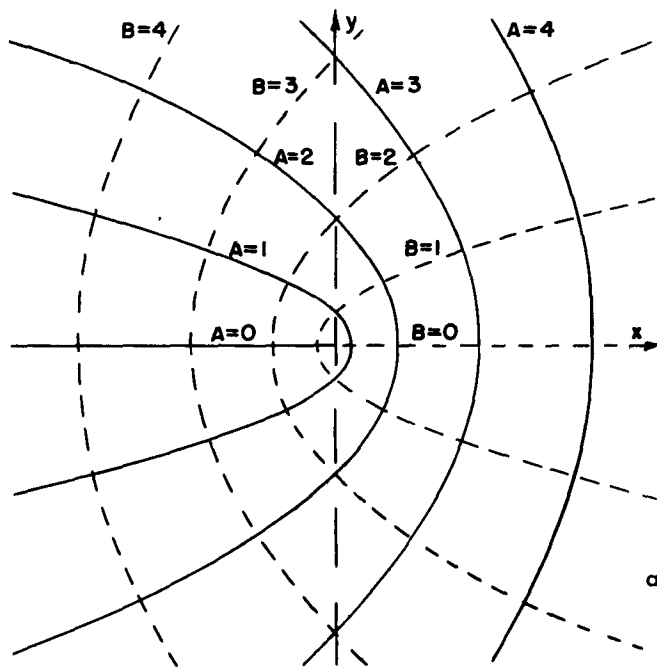


Fig. 7-4. Curves of constant $U(x, y)$ and $V(x, y)$ from Eq. (7.12).

The curves for $U(x, y) = A$ and $V(x, y) = B$, using U and V from Eq. (7.12), are plotted in Fig. 7-4. Again, there are many possible situations that could be described by these fields. One of the most interesting is the field near the edge of a thin plate. If the line $B = 0$ —to the right of the y -axis—represents a thin charged plate, the field lines near it are given by the curves for various values of A . The physical situation is shown in Fig. 7-5.

Further examples are

$$F(\bar{z}) = z^{3/2}, \quad (7.13)$$

which yields the field *outside* a rectangular corner

$$F(\bar{z}) = \log \bar{z}, \quad (7.14)$$

which yields the field for a line charge, and

$$F(\bar{z}) = 1/\bar{z}, \quad (7.15)$$

which gives the field for the two-dimensional analog of an electric dipole, i.e., two parallel line charges with opposite polarities, very close together.

We will not pursue this subject further in this course, but should emphasize that although the complex variable technique is often powerful, it is limited to two-dimensional problems; and also, it is an indirect method.

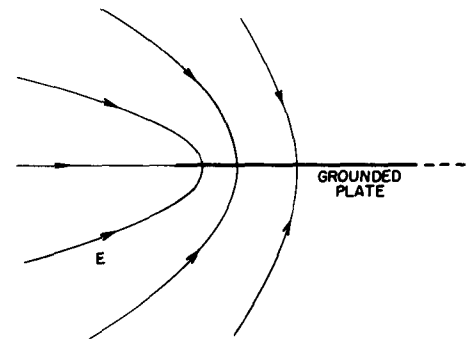


Fig. 7-5. The electric field near the edge of a thin grounded plate.

7-3 Plasma oscillations

We consider now some physical situations in which the field is determined neither by fixed charges nor by charges on conducting surfaces, but by a combination of two physical phenomena. In other words, the field will be governed simultaneously by two sets of equations: (1) the equations from electrostatics relating electric fields to charge distribution, and (2) an equation from another part of physics that determines the positions or motions of the charges in the presence of the field.

The first example that we will discuss is a dynamic one in which the motion of the charges is governed by Newton's laws. A simple example of such a situation occurs in a plasma, which is an ionized gas consisting of ions and free electrons distributed over a region in space. The ionosphere—an upper layer of the atmosphere—is an example of such a plasma. The ultraviolet rays from the sun knock

electrons off the molecules of the air, creating free electrons and ions. In such a plasma the positive ions are very much heavier than the electrons, so we may neglect the ionic motion, in comparison to that of the electrons.

Let n_0 be the density of electrons in the undisturbed, equilibrium state. This must also be the density of positive ions, since the plasma is electrically neutral (when undisturbed). Now we suppose that the electrons are somehow moved from equilibrium and ask what happens. If the density of the electrons in one region is increased, they will repel each other and tend to return to their equilibrium positions. As the electrons move toward their original positions they pick up kinetic energy, and instead of coming to rest in their equilibrium configuration, they overshoot the mark. They will oscillate back and forth. The situation is similar to what occurs in sound waves, in which the restoring force is the gas pressure. In a plasma, the restoring force is the electrical force on the electrons.

To simplify the discussion, we will worry only about a situation in which the motions are all in one dimension, say x . Let us suppose that the electrons originally at x are, at the instant t , displaced from their equilibrium positions by a small amount $s(x, t)$. Since the electrons have been displaced, their density will, in general, be changed. The change in density is easily calculated. Referring to Fig. 7-6, the electrons initially contained between the two planes a and b have moved and are now contained between the planes a' and b' . The number of electrons that were between a and b is proportional to $n_0 \Delta x$; the same number are now contained in the space whose width is $\Delta x + \Delta s$. The density has changed to

$$n = \frac{n_0 \Delta x}{\Delta x + \Delta s} = \frac{n_0}{1 + (\Delta s / \Delta x)}. \quad (7.16)$$

If the change in density is small, we can write [using the binomial expansion for $(1 + \epsilon)^{-1}$]

$$n = n_0 \left(1 - \frac{\Delta s}{\Delta x} \right). \quad (7.17)$$

We assume that the positive ions do not move appreciably (because of the much larger inertia), so their density remains n_0 . Each electron carries the charge $-q_e$, so the average charge density at any point is given by

$$\rho = -(n - n_0)q_e,$$

or

$$\rho = n_0 q_e \frac{ds}{dx} \quad (7.18)$$

(where we have written the differential form for $\Delta s / \Delta x$).

The charge density is related to the electric field by Maxwell's equations, in particular,

$$\nabla \cdot \mathbf{E} = \frac{\rho}{\epsilon_0}. \quad (7.19)$$

If the problem is indeed one-dimensional (and if there are no other fields but the one due to the displacements of the electrons), the electric field \mathbf{E} has a single component E_x . Equation (7.19), together with (7.18), gives

$$\frac{\partial E_x}{\partial x} = \frac{n_0 q_e}{\epsilon_0} \frac{\partial s}{\partial x}. \quad (7.20)$$

Integrating Eq. (7.20) gives

$$E_x = \frac{n_0 q_e}{\epsilon_0} s + K. \quad (7.21)$$

Since $E_x = 0$ when $s = 0$, the integration constant K is zero.

The force on an electron in the displaced position is

$$F_x = -\frac{n_0 q_e^2}{\epsilon_0} s, \quad (7.22)$$

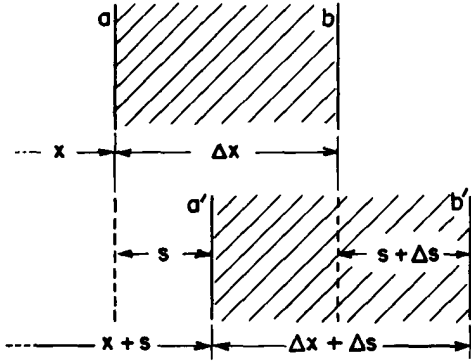


Fig. 7-6. Motion in a plasma wave. The electrons at the plane a move to a' , and those at b move to b' .

a restoring force proportional to the displacement s of the electron. This leads to a harmonic oscillation of the electrons. The equation of motion of a displaced electron is

$$m_e \frac{d^2 s}{dt^2} = -\frac{n_0 q_e^2}{\epsilon_0} s. \quad (7.23)$$

We find that s will vary harmonically. Its time variation will be as $\cos \omega t$, or—using the exponential notation of Vol. I—as

$$e^{i\omega_p t}. \quad (7.24)$$

The frequency of oscillation ω_p is determined from (7.23):

$$\omega_p^2 = \frac{n_0 q_e^2}{\epsilon_0 m_e}, \quad (7.25)$$

and is called the *plasma frequency*. It is a characteristic number of the plasma.

When dealing with electron charges many people prefer to express their answers in terms of a quantity e^2 defined by

$$e^2 = \frac{q_e^2}{4\pi\epsilon_0} = 2.3068 \times 10^{-28} \text{ newton} \cdot \text{meter}^2. \quad (7.26)$$

Using this convention, Eq. (7.25) becomes

$$\omega_p^2 = \frac{4\pi e^2 n_0}{m_e}, \quad (7.27)$$

which is the form you will find in most books.

Thus we have found that a disturbance of a plasma will set up free oscillations of the electrons about their equilibrium positions at the natural frequency ω_p , which is proportional to the square root of the density of the electrons. The plasma electrons behave like a resonant system, such as those we described in Chapter 23 of Vol. I.

This natural resonance of a plasma has some interesting effects. For example, if one tries to propagate a radiowave through the ionosphere, one finds that it can penetrate only if its frequency is higher than the plasma frequency. Otherwise the signal is reflected back. We must use high frequencies if we wish to communicate with a satellite in space. On the other hand, if we wish to communicate with a radio station beyond the horizon, we must use frequencies lower than the plasma frequency, so that the signal will be reflected back to the earth.

Another interesting example of plasma oscillations occurs in metals. In a metal we have a contained plasma of positive ions, and free electrons. The density n_0 is very high, so ω_p is also. But it should still be possible to observe the electron oscillations. Now, according to quantum mechanics, a harmonic oscillator with a natural frequency ω_p has energy levels which are separated by the energy increment $\hbar\omega_p$. If, then, one shoots electrons through, say, an aluminum foil, and makes very careful measurements of the electron energies on the other side, one might expect to find that the electrons sometimes lose the energy $\hbar\omega_p$ to the plasma oscillations. This does indeed happen. It was first observed experimentally in 1936 that electrons with energies of a few hundred to a few thousand electron volts lost energy in jumps when scattering from or going through a thin metal foil. The effect was not understood until 1953 when Bohm and Pines* showed that the observations could be explained in terms of quantum excitations of the plasma oscillations in the metal.

* For some recent work and a bibliography see C. J. Powell and J. B. Swann, *Phys. Rev.* **115**, 869 (1959).

7-4 Colloidal particles in an electrolyte

We turn to another phenomenon in which the locations of charges is governed by a potential that arises in part from the same charges. The resulting effects influence in an important way the behavior of colloids. A colloid consists of a suspension in water of small charged particles which, though microscopic, from an atomic point of view are still very large. If the colloidal particles were not charged, they would tend to coagulate into large lumps; but because of their charge, they repel each other and remain in suspension.

Now if there is also some salt dissolved in the water, it will be dissociated into positive and negative ions. (Such a solution of ions is called an electrolyte.) The negative ions are attracted to the colloid particles (assuming their charge is positive) and the positive ions are repelled. We will determine how the ions which surround such a colloidal particle are distributed in space.

To keep the ideas simple, we will again solve only a one-dimensional case. If we think of a colloidal particle as a sphere having a very large radius—on an atomic scale!—we can then treat a small part of its surface as a plane. (Whenever one is trying to understand a new phenomenon it is a good idea to take a somewhat oversimplified model; then, having understood the problem with that model, one is better able to proceed to tackle the more exact calculation.)

We suppose that the distribution of ions generates a charge density $\rho(x)$, and an electrical potential ϕ , related by the electrostatic law $\nabla^2\phi = -\rho/\epsilon_0$ or, for fields that vary in only one dimension, by

$$\frac{d^2\phi}{dx^2} = -\frac{\rho}{\epsilon_0}. \quad (7.28)$$

Now supposing there were such a potential $\phi(x)$, how would the ions distribute themselves in it? This we can determine by the principles of statistical mechanics. Our problem then is to determine ϕ so that the resulting charge density from statistical mechanics *also* satisfies (7.28).

According to statistical mechanics (see Chapter 40, Vol. I), particles in thermal equilibrium in a force field are distributed in such a way that the density n of particles at the position x is given by

$$n(x) = n_0 e^{-U(x)/kT}, \quad (7.29)$$

where $U(x)$ is the potential energy, k is Boltzmann's constant, and T is the absolute temperature.

We assume that the ions carry one electronic charge, positive or negative. At the distance x from the surface of a colloidal particle, a positive ion will have potential energy $q_e\phi(x)$, so that

$$U(x) = q_e\phi(x).$$

The density of positive ions, n_+ , is then

$$n_+(x) = n_0 e^{-q_e\phi(x)/kT}.$$

Similarly, the density of negative ions is

$$n_-(x) = n_0 e^{+q_e\phi(x)/kT}.$$

The total charge density is

$$\rho = q_e n_+ - q_e n_-,$$

or

$$\rho = q_e n_0 (e^{-q_e\phi/kT} - e^{+q_e\phi/kT}). \quad (7.30)$$

Combining this with Eq. (7.28), we find that the potential ϕ must satisfy

$$\frac{d^2\phi}{dx^2} = -\frac{q_e n_0}{\epsilon_0} (e^{-q_e\phi/kT} - e^{+q_e\phi/kT}). \quad (7.31)$$

This equation is readily solved in general [multiply both sides by $2(d\phi/dx)$, and integrate with respect to x], but to keep the problem as simple as possible, we will consider here only the limiting case in which the potentials are small or the temperature T is high. The case where ϕ is small corresponds to a dilute solution. For these cases the exponent is small, and we can approximate

$$e^{\pm q_e\phi/kT} = 1 \pm \frac{q_e\phi}{kT}. \quad (7.32)$$

Equation (7.31) then gives

$$\frac{d^2\phi}{dx^2} = + \frac{2n_0q_e^2}{\epsilon_0kT} \phi(x). \quad (7.33)$$

Notice that this time the sign on the right is positive. The solutions for ϕ are not oscillatory, but exponential.

The general solution of Eq. (7.33) is

$$\phi = Ae^{-x/D} + Be^{+x/D}, \quad (7.34)$$

with

$$D^2 = \frac{\epsilon_0kT}{2n_0q_e^2}. \quad (7.35)$$

The constants A and B must be determined from the conditions of the problem. In our case, B must be zero; otherwise the potential would go to infinity for large x . So we have that

$$\phi = Ae^{-x/D}, \quad (7.36)$$

in which A is the potential at $x = 0$, the surface of the colloidal particle.

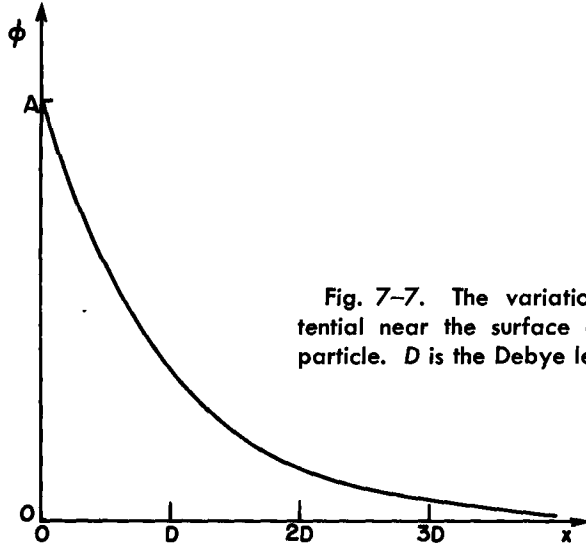


Fig. 7-7. The variation of the potential near the surface of a colloidal particle. D is the Debye length.

The potential decreases by a factor $1/e$ each time the distance increases by D , as shown in the graph of Fig. 7-7. The number D is called the *Debye length*, and is a measure of the thickness of the ion sheath that surrounds a large charged particle in an electrolyte. Equation (7.36) says that the sheath gets thinner with increasing concentration of the ions (n_0) or with decreasing temperature.

The constant A in Eq. (7.36) is easily obtained if we know the surface charge σ on the colloid particle. We know that

$$E_n = E_x(0) = \frac{\sigma}{\epsilon_0}. \quad (7.37)$$

But E is also the gradient of ϕ :

$$E_x(0) = - \left. \frac{\partial\phi}{\partial x} \right|_0 = + \frac{A}{D}, \quad (7.38)$$

from which we get

$$A = \frac{\sigma D}{\epsilon_0}. \quad (7.39)$$

Using this result in (7.36), we find (by taking $x = 0$) that the potential of the colloidal particle is

$$\phi(0) = \frac{\sigma D}{\epsilon_0}. \quad (7.40)$$

You will notice that this potential is the same as the potential difference across a condenser with a plate spacing D and a surface charge density σ .

We have said that the colloidal particles are kept apart by their electrical repulsion. But now we see that the field a little way from the surface of a particle is reduced by the ion sheath that collects around it. If the sheaths get thin enough, the particles have a good chance of knocking against each other. They will then stick, and the colloid will coagulate and precipitate out of the liquid. From our analysis, we understand why adding enough salt to a colloid should cause it to precipitate out. The process is called "salting out a colloid."

Another interesting example is the effect that a salt solution has on protein molecules. A protein molecule is a long, complicated, and flexible chain of amino acids. The molecule has various charges on it, and it sometimes happens that there is a net charge, say negative, which is distributed along the chain. Because of mutual repulsion of the negative charges, the protein chain is kept stretched out. Also, if there are other similar chain molecules present in the solution, they will be kept apart by the same repulsive effects. We can, therefore, have a suspension of chain molecules in a liquid. But if we add salt to the liquid we change the properties of the suspension. As salt is added to the solution, decreasing the Debye distance, the chain molecules can approach one another, and can also coil up. If enough salt is added to the solution, the chain molecules will precipitate out of the solution. There are many chemical effects of this kind that can be understood in terms of electrical forces.

7-5 The electrostatic field of a grid

As our last example, we would like to describe another interesting property of electric fields. It is one which is made use of in the design of electrical instruments, in the construction of vacuum tubes, and for other purposes. This is the character of the electric field near a grid of charged wires. To make the problem as simple as possible, let us consider an array of parallel wires lying in a plane, the wires being infinitely long and with a uniform spacing between them.

If we look at the field a large distance above the plane of the wires, we see a constant electric field, just as though the charge were uniformly spread over a plane. As we approach the grid of wires, the field begins to deviate from the uniform field we found at large distances from the grid. We would like to estimate how close to the grid we have to be in order to see appreciable variations in the potential. Figure 7-8 shows a rough sketch of the equipotential surfaces at various distances from the grid. The closer we get to the grid, the larger the variations. As we travel parallel to the grid, we observe that the field fluctuates in a periodic manner.

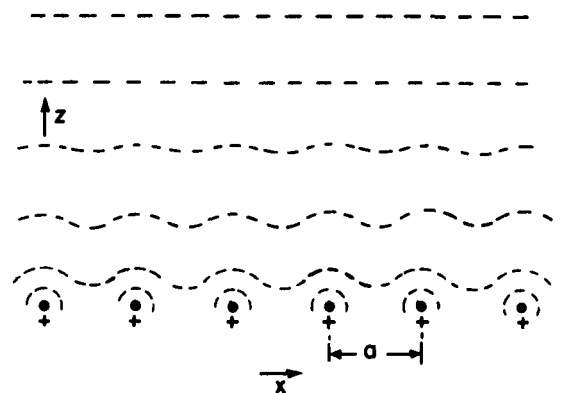


Fig. 7-8. Equipotential surfaces above a uniform grid of charged wires.

Now we have seen (Chapter 50, Vol. I) that any periodic quantity can be expressed as a sum of sine waves (Fourier's theorem). Let's see if we can find a suitable harmonic function that satisfies our field equations.

If the wires lie in the xy -plane and run parallel to the y -axis, then we might try terms like

$$\phi(x, z) = F_n(z) \cos \frac{2\pi nx}{a}, \quad (7.41)$$

where a is the spacing of the wires and n is the harmonic number. (We have assumed long wires, so there should be no variation with y .) A complete solution would be made up of a sum of such terms for $n = 1, 2, 3, \dots$

If this is to be a valid potential, it must satisfy Laplace's equation in the region above the wires (where there are no charges). That is,

$$\frac{\partial^2 \phi}{\partial x^2} + \frac{\partial^2 \phi}{\partial z^2} = 0.$$

Trying this equation on the ϕ in (7.41), we find that

$$-\frac{4\pi^2 n^2}{a^2} F_n(z) \cos \frac{2\pi nx}{a} + \frac{d^2 F_n}{dz^2} \cos \frac{2\pi nx}{a} = 0, \quad (7.42)$$

or that $F_n(z)$ must satisfy

$$\frac{d^2 F_n}{dz^2} = \frac{4\pi^2 n^2}{a^2} F_n. \quad (7.43)$$

So we must have

$$F_n = A_n e^{-z/z_0}, \quad (7.44)$$

where

$$z_0 = \frac{a}{2\pi n}. \quad (7.45)$$

We have found that if there is a Fourier component of the field of harmonic n , that component will decrease exponentially with a characteristic distance $z_0 = a/2\pi n$. For the first harmonic ($n = 1$), the amplitude falls by the factor $e^{-2\pi}$ (a large decrease) each time we increase z by one grid spacing a . The other harmonics fall off even more rapidly as we move away from the grid. We see that if we are only a few times the distance a away from the grid, the field is very nearly uniform, i.e., the oscillating terms are small. There would, of course, always remain the "zero harmonic" field

$$\phi_0 = -E_0 z$$

to give the uniform field at large z . For a complete solution, we would combine this term with a sum of terms like (7.41) with F_n from (7.44). The coefficients A_n would be adjusted so that the total sum would, when differentiated, give an electric field that would fit the charge density λ of the grid wires.

The method we have just developed can be used to explain why electrostatic shielding by means of a screen is often just as good as with a solid metal sheet. Except within a distance from the screen a few times the spacing of the screen wires, the fields inside a closed screen are zero. We see why copper screen—lighter and cheaper than copper sheet—is often used to shield sensitive electrical equipment from external disturbing fields.

Electrostatic Energy

8-1 The electrostatic energy of charges. A uniform sphere

In the study of mechanics, one of the most interesting and useful discoveries was the law of the conservation of energy. The expressions for the kinetic and potential energies of a mechanical system helped us to discover connections between the states of a system at two different times without having to look into the details of what was occurring in between. We wish now to consider the energy of electrostatic systems. In electricity also the principle of the conservation of energy will be useful for discovering a number of interesting things.

The law of the energy of interaction in electrostatics is very simple; we have, in fact, already discussed it. Suppose we have two charges q_1 and q_2 separated by the distance r_{12} . There is some energy in the system, because a certain amount of work was required to bring the charges together. We have already calculated the work done in bringing two charges together from a large distance. It is

$$\frac{q_1 q_2}{4\pi\epsilon_0 r_{12}}. \quad (8.1)$$

We also know, from the principle of superposition, that if we have many charges present, the total force on any charge is the sum of the forces from the others. It follows, therefore, that the total energy of a system of a number of charges is the sum of terms due to the mutual interaction of each pair of charges. If q_i and q_j are any two of the charges and r_{ij} is the distance between them (Fig. 8-1), the energy of that particular pair is

$$\frac{q_i q_j}{4\pi\epsilon_0 r_{ij}}. \quad (8.2)$$

The total electrostatic energy U is the sum of the energies of all possible pairs of charges:

$$U = \sum_{\text{all pairs}} \frac{q_i q_j}{4\pi\epsilon_0 r_{ij}}. \quad (8.3)$$

If we have a distribution of charge specified by a charge density ρ , the sum of Eq. (8.3) is, of course, to be replaced by an integral.

We shall concern ourselves with two aspects of this energy. One is the *application* of the concept of energy to electrostatic problems; the other is the *evaluation* of the energy in different ways. Sometimes it is easier to compute the work done for some special case than to evaluate the sum in Eq. (8.3), or the corresponding integral. As an example, let us calculate the energy required to assemble a sphere of charge with a uniform charge density. The energy is just the work done in gathering the charges together from infinity.

Imagine that we assemble the sphere by building up a succession of thin spherical layers of infinitesimal thickness. At each stage of the process, we gather a small amount of charge and put it in a thin layer from r to $r + dr$. We continue the process until we arrive at the final radius a (Fig. 8-2). If Q_r is the charge of the sphere when it has been built up to the radius r , the work done in bringing a charge dQ to it is

$$dU = \frac{Q_r dQ}{4\pi\epsilon_0 r}. \quad (8.4)$$

8-1

8-1 The electrostatic energy of charges. A uniform sphere

8-2 The energy of a condenser. Forces on charged conductors

8-3 The electrostatic energy of an ionic crystal

8-4 Electrostatic energy in nuclei

8-5 Energy in the electrostatic field

8-6 The energy of a point charge

Review: Chapter 4, Vol. I, Conservation of Energy
Chapters 13 and 14, Vol. I, Work and Potential Energy

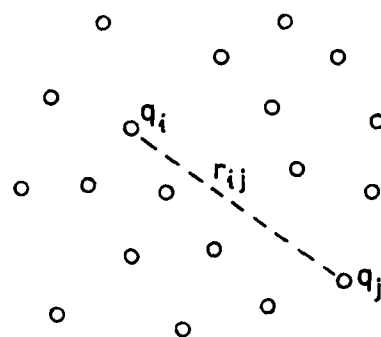


Fig. 8-1. The electrostatic energy of a system of particles is the sum of the electrostatic energy of each pair.

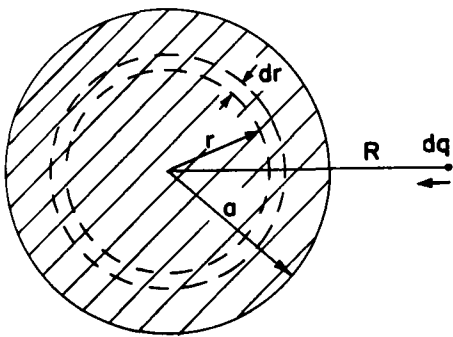


Fig. 8-2. The energy of a uniform sphere of charge can be computed by imagining that it is assembled from successive spherical shells.

If the density of charge in the sphere is ρ , the charge Q_r is

$$Q_r = \rho \cdot \frac{4}{3} \pi r^3,$$

and the charge dQ is

$$dQ = \rho \cdot 4\pi r^2 dr.$$

Equation (8.4) becomes

$$dU = \frac{4\pi\rho^2 r^4 dr}{3\epsilon_0}. \quad (8.5)$$

The total energy required to assemble the sphere is the integral of dU from $r = 0$ to $r = a$, or

$$U = \frac{4\pi\rho^2 a^5}{15\epsilon_0}. \quad (8.6)$$

Or if we wish to express the result in terms of the total charge Q of the sphere,

$$U = \frac{3}{5} \frac{Q^2}{4\pi\epsilon_0 a}. \quad (8.7)$$

The energy is proportional to the square of the total charge and inversely proportional to the radius. We can also interpret Eq. (8.7) as saying that the average of $(1/r_{ij})$ for all pairs of points in the sphere is $3/5a$.

8-2 The energy of a condenser. Forces on charged conductors

We consider now the energy required to charge a condenser. If the charge Q has been taken from one of the conductors of a condenser and placed on the other, the potential difference between them is

$$V = \frac{Q}{C}, \quad (8.8)$$

where C is the capacity of the condenser. How much work is done in charging the condenser? Proceeding as for the sphere, we imagine that the condenser has been charged by transferring charge from one plate to the other in small increments dQ . The work required to transfer the charge dQ is

$$dU = V dQ.$$

Taking V from Eq. (8.8), we write

$$dU = \frac{Q dQ}{C}.$$

Or integrating from zero charge to the final charge Q , we have

$$U = \frac{1}{2} \frac{Q^2}{C}. \quad (8.9)$$

This energy can also be written as

$$U = \frac{1}{2} CV^2. \quad (8.10)$$

Recalling that the capacity of a conducting sphere (relative to infinity) is

$$C_{\text{sphere}} = 4\pi\epsilon_0 a,$$

we can immediately get from Eq. (8.9) the energy of a charged sphere,

$$U = \frac{1}{2} \frac{Q^2}{4\pi\epsilon_0 a}. \quad (8.11)$$

This, of course, is also the energy of a thin *spherical shell* of total charge Q and is just 5/6 of the energy of a *uniformly charged sphere*, Eq. (8.7).

We now consider applications of the idea of electrostatic energy. Consider the following questions: What is the force between the plates of a condenser? Or what is the torque about some axis of a charged conductor in the presence of another with opposite charge? Such questions are easily answered by using our result Eq. (8.9) for electrostatic energy of a condenser, together with the principle of virtual work (Chapters 4, 13, and 14 of Vol. I).

Let's use this method for determining the force between the plates of a parallel-plate condenser. If we imagine that the spacing of the plates is increased by the small amount Δz , then the mechanical work done from the outside in moving the plates would be

$$\Delta W = F \Delta z, \quad (8.12)$$

where F is the force between the plates. This work must be equal to the change in the electrostatic energy of the condenser.

By Eq. (8.9), the energy of the condenser was originally

$$U = \frac{1}{2} \frac{Q^2}{C}.$$

The change in energy (if we do not let the charge change) is

$$\Delta U = \frac{1}{2} Q^2 \Delta \left(\frac{1}{C} \right). \quad (8.13)$$

Equating (8.12) and (8.13), we have

$$F \Delta z = \frac{Q^2}{2} \Delta \left(\frac{1}{C} \right). \quad (8.14)$$

This can also be written as

$$F \Delta z = - \frac{Q^2}{2C^2} \Delta C. \quad (8.15)$$

The force, of course, results from the attraction of the charges on the plates, but we see that we do not have to worry in detail about how they are distributed; everything we need is taken care of in the capacity C .

It is easy to see how the idea is extended to conductors of any shape, and for other components of the force. In Eq. (8.14), we replace F by the component we are looking for, and we replace Δz by a small displacement in the corresponding direction. Or if we have an electrode with a pivot and we want to know the torque τ , we write the virtual work as

$$\Delta W = \tau \Delta \theta,$$

where $\Delta \theta$ is a small angular displacement. Of course, $\Delta(1/C)$ must be the change in $1/C$ which corresponds to $\Delta \theta$. We could, in this way, find the torque on the movable plates in a variable condenser of the type shown in Fig. 8-3.

Returning to the special case of a parallel-plate condenser, we can use the formula we derived in Chapter 6 for the capacity:

$$\frac{1}{C} = \frac{d}{\epsilon_0 A}, \quad (8.16)$$

where A is the area of each plate. If we increase the separation by Δz ,

$$\Delta \left(\frac{1}{C} \right) = \frac{\Delta z}{\epsilon_0 A}.$$

From Eq. (8.14) we get that the force between the plates is

$$F = \frac{Q^2}{2\epsilon_0 A}. \quad (8.17)$$

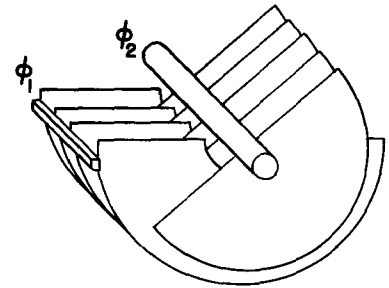


Fig. 8-3. What is the torque on a variable capacitor?

Let's look at Eq. (8.17) a little more closely and see if we can tell how the force arises. If for the charge on one plate we write

$$Q = \sigma A,$$

Eq. (8.17) can be rewritten as

$$F = \frac{1}{2} Q \frac{\sigma}{\epsilon_0}.$$

Or, since the electric field between the plates is

$$E_0 = \frac{\sigma}{\epsilon_0},$$

then

$$F = \frac{1}{2} Q E_0. \quad (8.18)$$

One would immediately guess that the force acting on one plate is the charge Q on the plate times the field acting on the charge. But we have a surprising factor of one-half. The reason is that E_0 is not the field *at* the charges. If we imagine that the charge at the surface of the plate occupies a thin layer, as indicated in Fig. 8-4, the field will vary from zero at the inner boundary of the layer to E_0 in the space outside of the plate. The average field acting on the surface charges is $E_0/2$. That is why the factor one-half is in Eq. (8.18).

You should notice that in computing the virtual work we have assumed that the charge on the condenser was constant—that it was not electrically connected to other objects, and so the total charge could not change.

Suppose we had imagined that the condenser was held at a constant potential difference as we made the virtual displacement. Then we should have taken

$$U = \frac{1}{2} C V^2$$

and in place of Eq. (8.15) we would have had

$$F \Delta z = \frac{1}{2} V^2 \Delta C,$$

which gives a force equal in magnitude to the one in Eq. (8.15) (because $V = Q/C$), but with the opposite sign! Surely the force between the condenser plates doesn't reverse in sign as we disconnect it from its charging source. Also, we know that two plates with opposite electrical charges must attract. The principle of virtual work has been incorrectly applied in the second case—we have not taken into account the virtual work done on the charging source. That is, to keep the potential constant at V as the capacity changes, a charge $V \Delta C$ must be supplied by a source of charge. But this charge is supplied at a potential V , so the work done by the electrical system which keeps the potential constant is $V^2 \Delta C$. The mechanical work $F \Delta z$ plus this electrical work $V^2 \Delta C$ together make up the change in the total energy $\frac{1}{2} V^2 \Delta C$ of the condenser. Therefore $F \Delta z$ is $-\frac{1}{2} V^2 \Delta C$, as before.

8-3 The electrostatic energy of an ionic crystal

We now consider an application of the concept of electrostatic energy in atomic physics. We cannot easily measure the forces between atoms, but we are often interested in the energy differences between one atomic arrangement and another, as, for example, the energy of a chemical change. Since atomic forces are basically electrical, chemical energies are in large part just electrostatic energies.

Let's consider, for example, the electrostatic energy of an ionic lattice. An ionic crystal like NaCl consists of positive and negative ions which can be thought of as rigid spheres. They attract electrically until they begin to touch; then there is a repulsive force which goes up very rapidly if we try to push them closer together.

For our first approximation, therefore, we imagine a set of rigid spheres that represent the atoms in a salt crystal. The structure of the lattice has been determined by x-ray diffraction. It is a cubic lattice—like a three-dimensional

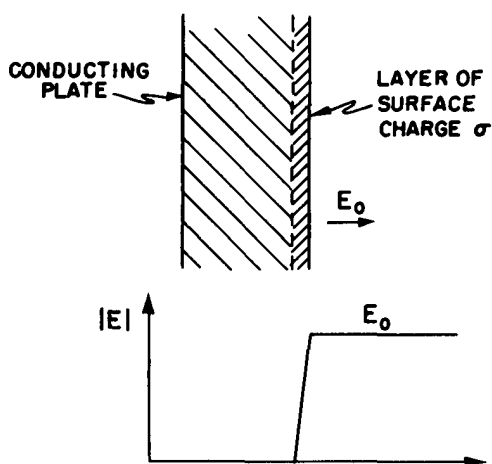


Fig. 8-4. The field at the surface of a conductor varies from zero to $E_0 = \sigma/\epsilon_0$, as one passes through the layer of surface charge.

checkerboard. Figure 8-5 shows a cross-sectional view. The spacing of the ions is 2.81 \AA ($=2.81 \times 10^{-8} \text{ cm}$).

If our picture of this system is correct, we should be able to check it by asking the following question: How much energy will it take to pull all these ions apart—that is, to separate the crystal completely into ions? This energy should be equal to the heat of vaporization of NaCl plus the energy required to dissociate the molecules into ions. This total energy to separate NaCl to ions is determined experimentally to be 7.92 electron volts per molecule. Using the conversion

$$1 \text{ ev} = 1.602 \times 10^{-19} \text{ joule},$$

and Avogadro's number for the number of molecules in a mole,

$$N_0 = 6.02 \times 10^{23},$$

the energy of vaporization can also be given as

$$W = 7.64 \times 10^5 \text{ joules/mole}.$$

Physical chemists prefer for an energy unit the kilocalorie, which is 4190 joules; so that 1 ev per molecule is 23 kilocalories per mole. A chemist would then say that the dissociation energy of NaCl is

$$W = 183 \text{ kcal/mole}.$$

Can we obtain this chemical energy theoretically by computing how much work it would take to pull apart the crystal? According to our theory, this work is the sum of the potential energies of all the pairs of ions. The easiest way to figure out this sum is to pick out a particular ion and compute its potential energy with each of the other ions. That will give us *twice* the energy per ion, because the energy belongs to the *pairs* of charges. If we want the energy to be associated with one particular ion, we should take half the sum. But we really want the energy *per molecule*, which contains two ions, so that the sum we compute will give directly the energy per molecule.

The energy of an ion with one of its nearest neighbors is e^2/a , where $e^2 = q_e^2/4\pi\epsilon_0$ and a is the center-to-center spacing between ions. (We are considering monovalent ions.) This energy is 5.12 ev, which we already see is going to give us a result of the correct order of magnitude. But it is still a long way from the infinite sum of terms we need.

Let's begin by summing all the terms from the ions along a straight line. Considering that the ion marked Na in Fig. 8-5 is our special ion, we shall consider first those ions on a horizontal line with it. There are two nearest Cl ions with negative charges, each at the distance a . Then there are two positive ions at the distance $2a$, etc. Calling the energy of this sum U_1 , we write

$$\begin{aligned} U_1 &= \frac{e^2}{a} \left(-\frac{2}{1} + \frac{2}{2} - \frac{2}{3} + \frac{2}{4} + \cdots \right) \\ &= -\frac{2e^2}{a} \left(1 - \frac{1}{2} + \frac{1}{3} - \frac{1}{4} + \cdots \right). \end{aligned} \quad (8.19)$$

The series converges slowly, so it is difficult to evaluate numerically, but it is known to be equal to $\ln 2$. So

$$U_1 = -\frac{2e^2}{a} \ln 2 = -1.386 \frac{e^2}{a}. \quad (8.20)$$

Now consider the next adjacent line of ions above. The nearest is negative and at the distance a . Then there are two positives at the distance $\sqrt{2}a$. The next pair are at the distance $\sqrt{5}a$, the next at $\sqrt{10}a$, and so on. So for the whole line we get the series

$$\frac{e^2}{a} \left(-\frac{1}{1} + \frac{2}{\sqrt{2}} - \frac{2}{\sqrt{5}} + \frac{2}{\sqrt{10}} \cdots \right). \quad (8.21)$$

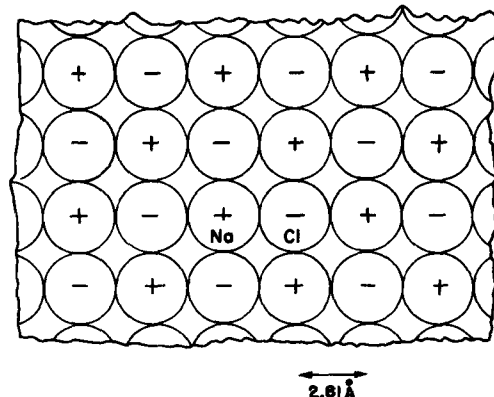


Fig. 8-5. Cross section of a salt crystal on an atomic scale. The checkerboard arrangement of Na and Cl ions is the same in the two cross sections perpendicular to the one shown. (See Vol. I, Fig. 1-7.)

There are *four* such lines: above, below, in front, and in back. Then there are the four lines which are the nearest lines on diagonals, and on and on.

If you work patiently through for all the lines, and then take the sum, you find that the grand total is

$$U = 1.747 \frac{e^2}{a},$$

which is just somewhat more than what we obtained in (8.20) for the first line. Using $e^2/a = 5.12$ ev, we get

$$U = 8.94 \text{ ev.}$$

Our answer is about 10% above the experimentally observed energy. It shows that our idea that the whole lattice is held together by electrical Coulomb forces is fundamentally correct. This is the first time that we have obtained a specific property of a macroscopic substance from a knowledge of atomic physics. We will do much more later. The subject that tries to understand the behavior of bulk matter in terms of the laws of atomic behavior is called *solid-state physics*.

Now what about the error in our calculation? Why is it not exactly right? It is because of the repulsion between the ions at close distances. They are not perfectly rigid spheres, so when they are close together they are partly squashed. They are not very soft, so they squash only a little bit. Some energy, however, is used in deforming them, and when the ions are pulled apart this energy is released. The actual energy needed to pull the ions apart is a little less than the energy that we calculated; the repulsion helps in overcoming the electrostatic attraction.

Is there any way we can make an allowance for this contribution? We could if we knew the law of the repulsive force. We are not ready to analyze the details of this repulsive mechanism, but we can get some idea of its characteristics from some large-scale measurements. From a measurement of the *compressibility* of the whole crystal, it is possible to obtain a quantitative idea of the law of repulsion between the ions and therefore of its contribution to the energy. In this way it has been found that this contribution must be 1/9.4 of the contribution from the electrostatic attraction and, of course, of opposite sign. If we subtract this contribution from the pure electrostatic energy, we obtain 7.99 ev for the dissociation energy per molecule. It is much closer to the observed result of 7.92 ev, but still not in perfect agreement. There is one more thing we haven't taken into account: we have made no allowance for the kinetic energy of the crystal vibrations. If a correction is made for this effect, very good agreement with the experimental number is obtained. The ideas are then correct; the major contribution to the energy of a crystal like NaCl is electrostatic.

8-4 Electrostatic energy in nuclei

We will now take up another example of electrostatic energy in atomic physics, the electrical energy of atomic nuclei. Before we do this we will have to discuss some properties of the main forces (called nuclear forces) that hold the protons and neutrons together in a nucleus. In the early days of the discovery of nuclei—and of the neutrons and protons that make them up—it was hoped that the law of the strong, nonelectrical part of the force between, say, a proton and another proton would have some simple law, like the inverse square law of electricity. For once one had determined this law of force, and the corresponding ones between a proton and a neutron, and a neutron and a neutron, it would be possible to describe theoretically the complete behavior of these particles in nuclei. Therefore a big program was started for the study of the scattering of protons, in the hope of finding the law of force between them; but after thirty years of effort, nothing simple has emerged. A considerable knowledge of the force between proton and proton has been accumulated, but we find that the force is as complicated as it can possibly be.

What we mean by “as complicated as it can be” is that the force depends on as many things as it possibly can.

First, the force is not a simple function of the distance between the two protons. At large distances there is an attraction, but at closer distances there is a repulsion. The distance dependence is a complicated function, still imperfectly known.

Second, the force depends on the orientation of the protons' spin. The protons have a spin, and any two interacting protons may be spinning with their angular momenta in the same direction or in opposite directions. And the force is different when the spins are parallel from what it is when they are antiparallel, as in (a) and (b) of Fig. 8-6. The difference is quite large; it is not a small effect.

Third, the force is considerably different when the separation of the two protons is in the direction *parallel* to their spins, as in (c) and (d) of Fig. 8-6, than it is when the separation is in a direction *perpendicular* to the spins, as in (a) and (b).

Fourth, the force depends, as it does in magnetism, on the velocity of the protons, only much more strongly than in magnetism. And this velocity-dependent force is not a relativistic effect; it is strong even at speeds much less than the speed of light. Furthermore, this part of the force depends on other things besides the magnitude of the velocity. For instance, when a proton is moving near another proton, the force is different when the orbital motion has the same direction of rotation as the spin, as in (e) of Fig. 8-6, than when it has the opposite direction of rotation, as in (f). This is called the "spin orbit" part of the force.

The force between a proton and a neutron and between a neutron and a neutron are also equally complicated. To this day we do not know the machinery behind these forces—that is to say, any simple way of understanding them.

There is, however, one important way in which the nucleon forces are *simpler* than they could be. That is that the *nuclear* force between two neutrons is the same as the force between a proton and a neutron, which is the same as the force between two protons! If, in any nuclear situation, we replace a proton by a neutron (or vice versa), the *nuclear interactions* are not changed. The "fundamental reason" for this equality is not known, but it is an example of an important principle that can be extended also to the interaction laws of other strongly interacting particles—such as the π -mesons and the "strange" particles.

This fact is nicely illustrated by the locations of the energy levels in similar nuclei. Consider a nucleus like B^{11} (boron-eleven), which is composed of five protons and six neutrons. In the nucleus the eleven particles interact with one another in a most complicated dance. Now, there is one configuration of all the possible interactions which has the lowest possible energy; this is the normal state of the nucleus, and is called the *ground state*. If the nucleus is disturbed (for example, by being struck by a high-energy proton or other particle) it can be put into any number of other configurations, called *excited states*, each of which will have a characteristic energy that is higher than that of the ground state. In nuclear physics research, such as is carried on with Van de Graaff generator (for example, in Caltech's Kellogg and Sloan Laboratories), the energies and other properties of these excited states are determined by experiment. The energies of the fifteen lowest known excited states of B^{11} are shown in a one-dimensional graph on the left half of Fig. 8-7. The lowest horizontal line represents the ground state. The first excited state has an energy 2.14 Mev higher than the ground state, the next an energy 4.46 Mev higher than the ground state, and so on. The study of nuclear physics attempts to find an explanation for this rather complicated pattern of energies; there is as yet, however, no complete general theory of such nuclear energy levels.

If we replace one of the neutrons in B^{11} with a proton, we have the nucleus of an isotope of carbon, C^{11} . The energies of the lowest sixteen excited states of C^{11} have also been measured; they are shown in the right half of Fig. 8-7. (The broken lines indicate levels for which the experimental information is questionable.)

Looking at Fig. 8-7, we see a striking similarity between the pattern of the energy levels in the two nuclei. The first excited states are about 2 Mev above the ground states. There is a large gap of about 2.3 Mev to the second excited state, then a small jump of only 0.5 Mev to the third level. Again, between the fourth and fifth levels, a big jump; but between the fifth and sixth a tiny separation of the

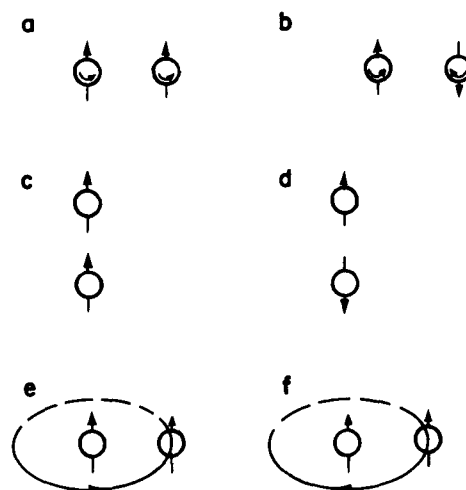


Fig. 8-6. The force between two protons depends on every possible parameter.

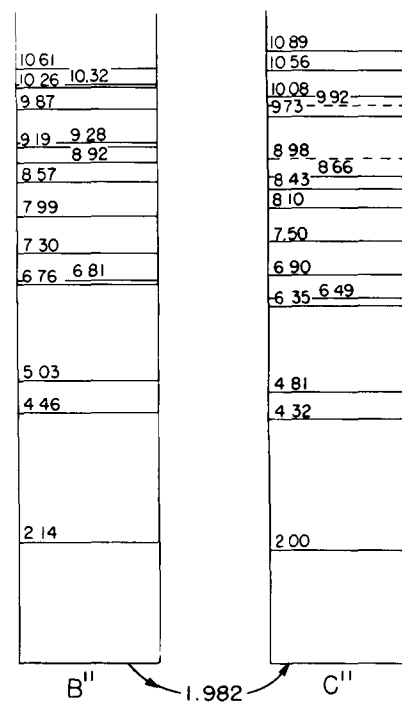


Fig. 8-7. The energy levels of B^{11} and C^{11} (energies in Mev). The ground state of C^{11} is 1.982 Mev higher than that of B^{11} .

order of 0.1 Mev. And so on. After about the tenth level, the correspondence seems to become lost, but can still be seen if the levels are labeled with their other defining characteristics—for instance, their angular momentum and what they do to lose their extra energy.

The striking similarity of the pattern of the energy levels of B^{11} and C^{11} is surely not just a coincidence. It must reveal some physical law. It shows, in fact, that even in the complicated situation in a nucleus, replacing a neutron by a proton makes very little change. This can mean only that the neutron-neutron and proton-proton forces must be nearly identical. Only then would we expect the nuclear configurations with five protons and six neutrons to be the same as with six protons and five neutrons.

Notice that the properties of these two nuclei tell us nothing about the neutron-proton force; there are the same number of neutron-proton combinations in both nuclei. But if we compare two other nuclei, such as C^{14} , which has six protons and eight neutrons, with N^{14} , which has seven of each, we find a similar correspondence of energy levels. So we can conclude that the p-p, n-n, and p-n forces are identical in all their complexities. There is an unexpected principle in the laws of nuclear forces. Even though the force between each pair of nuclear particles is very complicated, the force between the three possible different pairs is the same.

But there are some small differences. The levels do not correspond exactly; also, the ground state of C^{11} has an absolute energy (its mass) which is higher than the ground state of B^{11} by 1.982 Mev. All the other levels are also higher in absolute energy by this same amount. So the forces are not exactly equal. But we know very well that the *complete* forces are not exactly equal; there is an *electrical* force between two protons because each has a positive charge, while between two neutrons there is no such electrical force. Can we perhaps explain the differences between B^{11} and C^{11} by the fact that the electrical interaction of the protons is different in the two cases? Perhaps even the remaining minor differences in the levels are caused by electrical effects? Since the nuclear forces are so much stronger than the electrical force, electrical effects would have only a small perturbing effect on the energies of the levels.

In order to check this idea, or rather to find out what the consequences of this idea are, we first consider the difference in the ground-state energies of the two nuclei. To take a very simple model, we suppose that the nuclei are spheres of radius r (to be determined), containing Z protons. If we consider that a nucleus is like a sphere with uniform charge density, we would expect the electrostatic energy (from Eq. 8.7) to be

$$U = \frac{3}{5} \frac{(Zq_e)^2}{4\pi\epsilon_0 r}, \quad (8.22)$$

where q_e is the elementary charge of the proton. Since Z is five for B^{11} and six for C^{11} , their electrostatic energies would be different.

With such a small number of protons, however, Eq. (8.22) is not quite correct. If we compute the electrical energy between all pairs of protons, considered as points which we assume to be nearly uniformly distributed throughout the sphere, we find that in Eq. (8.22) the quantity Z^2 should be replaced by $Z(Z - 1)$, so the energy is

$$U = \frac{3}{5} \frac{Z(Z - 1)q_e^2}{4\pi\epsilon_0 r} = \frac{3}{5} \frac{Z(Z - 1)e^2}{r}. \quad (8.23)$$

If we knew the nuclear radius r , we could use (8.23) to find the electrostatic energy difference between B^{11} and C^{11} . But let's do the opposite; let's instead use the observed energy difference to compute the radius, assuming that the energy difference is all electrostatic in origin.

That is, however, not quite right. The energy difference of 1.982 Mev between the ground states of B^{11} and C^{11} includes the rest energies—that is, the energy mc^2 —of all the particles. In going from B^{11} to C^{11} , we replace a neutron by a proton, which has less mass. So part of the energy difference is the difference in the rest energies of a neutron and a proton, which is 0.784 Mev. The difference,

to be accounted for by electrostatic energy, is thus more than 1.982 Mev; it is

$$1.982 + 0.784 = 2.786 \text{ Mev.}$$

Using this energy in Eq. (8.23), for the radius of either B^{11} or C^{11} we find

$$r = 3.12 \times 10^{-13} \text{ cm.} \quad (8.24)$$

Does this number have any meaning? To see whether it does, we should compare it with some other determination of the radius of these nuclei. For example, we can make another measurement of the radius of a nucleus by seeing how it scatters fast particles. From such measurements it has been found, in fact, that the *density* of matter in all nuclei is nearly the same, i.e., their volumes are proportional to the number of particles they contain. If we let A be the number of protons and neutrons in a nucleus (a number very nearly proportional to its mass), it is found that its radius is given by

$$r = A^{1/3} r_0, \quad (8.25)$$

where

$$r_0 = 1.2 \times 10^{-13} \text{ cm.} \quad (8.26)$$

From these measurements we find that the radius of a B^{11} (or a C^{11}) nucleus is expected to be

$$r = (1.2 \times 10^{-13})(11)^{1/3} = 2.7 \times 10^{-13} \text{ cm.}$$

Comparing this result with (8.24), we see that our assumptions that the energy difference between B^{11} and C^{11} is electrostatic is fairly good; the discrepancy is only about 15% (not bad for our first nuclear computation!).

The reason for the discrepancy is probably the following. According to the current understanding of nuclei, an even number of nuclear particles—in the case of B^{11} , five neutrons together with five protons—makes a kind of *core*; when one more particle is added to this core, it revolves around on the outside to make a new spherical nucleus, rather than being absorbed. If this is so, we should have taken a different electrostatic energy for the additional proton. We should have taken the excess energy of C^{11} over B^{11} to be just

$$\frac{Z_B q_e^2}{4\pi\epsilon_0 a},$$

which is the energy needed to add one more proton to the outside of the core. This number is just 5/6 of what Eq. (8.23) predicts, so the new prediction for the radius is 5/6 of (8.24), which is in much closer agreement with what is directly measured.

We can draw two conclusions from this agreement. One is that the electrical laws appear to be working at dimensions as small as 10^{-13} cm. The other is that we have verified the remarkable coincidence that the nonelectrical part of the forces between proton and proton, neutron and neutron, and proton and neutron are all equal.

8-5 Energy in the electrostatic field

We now consider other methods of calculating electrostatic energy. They can all be derived from the basic relation Eq. (8.3), the sum, over all pairs of charges, of the mutual energies of each charge-pair. First we wish to write an expression for the energy of a charge distribution. As usual, we consider that each volume element dV contains the element of charge ρdV . Then Eq. (8.3) should be written

$$U = \frac{1}{2} \int_{\text{all space}} \frac{\rho(1)\rho(2)}{4\pi\epsilon_0 r_{12}} dV_1 dV_2. \quad (8.27)$$

Notice the factor $\frac{1}{2}$, which is introduced because in the double integral over dV_1 and dV_2 we have counted all pairs of charge elements twice. (There is no convenient way of writing an integral that keeps track of the pairs so that each pair is counted only once.) Next we notice that the integral over dV_2 in (8.27) is just the potential at (1). That is,

$$\int \frac{\rho(2)}{4\pi\epsilon_0 r_{12}} dV_2 = \phi(1),$$

so that (8.27) can be written as

$$U = \frac{1}{2} \int \rho(1)\phi(1) dV_1.$$

Or, since the point (2) no longer appears, we can simply write

$$U = \frac{1}{2} \int \rho\phi dV. \quad (8.28)$$

This equation can be interpreted as follows. The potential energy of the charge ρdV is the product of this charge and the potential at the same point. The total energy is therefore the integral over $\phi\rho dV$. But there is again the factor $\frac{1}{2}$. It is still required because we are counting energies twice. The mutual energies of two charges is the charge of one times the potential at it due to the other. Or, it can be taken as the second charge times the potential at it from the first. Thus for two point charges we would write

$$U = q_1\phi(1) = q_1 \frac{q_2}{4\pi\epsilon_0 r_{12}}$$

or

$$U = q_2\phi(2) = q_2 \frac{q_1}{4\pi\epsilon_0 r_{12}}.$$

Notice that we could also write

$$U = \frac{1}{2}[q_1\phi(1) + q_2\phi(2)]. \quad (8.29)$$

The integral in (8.28) corresponds to the sum of both terms in the brackets of (8.29). That is why we need the factor $\frac{1}{2}$.

An interesting question is: Where is the electrostatic energy located? One might also ask: Who cares? What is the meaning of such a question? If there is a pair of interacting charges, the combination has a certain energy. Do we need to say that the energy is located at one of the charges or the other, or at both, or in between? These questions may not make sense because we really know only that the total energy is conserved. The idea that the energy is located *somewhere* is not necessary.

Yet suppose that it *did* make sense to say, in general, that energy is located at a certain place, as it does for heat energy. We might then *extend* our principle of the conservation of energy with the idea that if the energy in a given volume changes, we should be able to account for the change by the flow of energy into or out of that volume. You realize that our early statement of the principle of the conservation of energy is still perfectly all right if some energy disappears at one place and appears somewhere else far away without anything passing (that is, without any special phenomena occurring) in the space between. We are, therefore, now discussing an extension of the idea of the conservation of energy. We might call it a principle of the *local* conservation of energy. Such a principle would say that the energy in any given volume changes only by the amount that flows into or out of the volume. It is indeed possible that energy is conserved locally in such a way. If it is, we would have a much more detailed law than the simple statement of the conservation of total energy. It does turn out that in nature *energy is conserved locally*. We can find formulas for where the energy is located and how it travels from place to place.

There is also a *physical* reason why it is imperative that we be able to say where energy is located. According to the theory of gravitation, all mass is a source

of gravitational attraction. We also know, by $E = mc^2$, that mass and energy are equivalent. All energy is, therefore, a source of gravitational force. If we could not locate the energy, we could not locate all the mass. We would not be able to say where the sources of the gravitational field are located. The theory of gravitation would be incomplete.

If we restrict ourselves to electrostatics there is really no way to tell where the energy is located. The complete Maxwell equations of electrodynamics give us much more information (although even then the answer is, strictly speaking, not unique.) We will therefore discuss this question in detail again in a later chapter. We will give you now only the result for the particular case of electrostatics. The energy is located in space, where the electric field is. This seems reasonable because we know that when charges are accelerated they radiate electric fields. We would like to say that when light or radiowaves travel from one point to another, they carry their energy with them. But there are no charges in the waves. So we would like to locate the energy where the electromagnetic field is and not at the charges from which it came. We thus describe the energy, not in terms of the charges, but in terms of the fields they produce. We can, in fact, show that Eq. (8.28) is *numerically* equal to

$$U = \frac{\epsilon_0}{2} \int \mathbf{E} \cdot \mathbf{E} dV. \quad (8.30)$$

We can then interpret this formula as saying that when an electric field is present, there is located in space an energy whose *density* (energy per unit volume) is

$$u = \frac{\epsilon_0}{2} \mathbf{E} \cdot \mathbf{E} = \frac{\epsilon_0 E^2}{2}. \quad (8.31)$$

This idea is illustrated in Fig. 8-8.

To show that Eq. (8.30) is consistent with our laws of electrostatics, we begin by introducing into Eq. (8.28) the relation between ρ and ϕ that we obtained in Chapter 6:

$$\rho = -\epsilon_0 \nabla^2 \phi.$$

We get

$$U = -\frac{\epsilon_0}{2} \int \phi \nabla^2 \phi dV. \quad (8.32)$$

Writing out the components of the integrand, we see that

$$\begin{aligned} \phi \nabla^2 \phi &= \phi \left(\frac{\partial^2 \phi}{\partial x^2} + \frac{\partial^2 \phi}{\partial y^2} + \frac{\partial^2 \phi}{\partial z^2} \right) \\ &= \frac{\partial}{\partial x} \left(\phi \frac{\partial \phi}{\partial x} \right) - \left(\frac{\partial \phi}{\partial x} \right)^2 + \frac{\partial}{\partial y} \left(\phi \frac{\partial \phi}{\partial y} \right) - \left(\frac{\partial \phi}{\partial y} \right)^2 + \frac{\partial}{\partial z} \left(\phi \frac{\partial \phi}{\partial z} \right) - \left(\frac{\partial \phi}{\partial z} \right)^2 \\ &= \nabla \cdot (\phi \nabla \phi) - (\nabla \phi) \cdot (\nabla \phi). \end{aligned} \quad (8.33)$$

Our energy integral is then

$$U = \frac{\epsilon_0}{2} \int (\nabla \phi) \cdot (\nabla \phi) dV - \frac{\epsilon_0}{2} \int \nabla \cdot (\phi \nabla \phi) dV.$$

We can use Gauss' theorem to change the second integral into a surface integral:

$$\int_{\text{vol.}} \nabla \cdot (\phi \nabla \phi) dV = \int_{\text{surface}} (\phi \nabla \phi) \cdot \mathbf{n} da. \quad (8.34)$$

We evaluate the surface integral in the case that the surface goes to infinity (so the volume integrals become integrals over all space), supposing that all the charges are located within some finite distance. The simple way to proceed is to take a spherical surface of enormous radius R whose center is at the origin of coordinates. We know that when we are very far away from all charges, ϕ varies as $1/R$ and $\nabla \phi$ as $1/R^2$. (Both will decrease even faster with R if there the net

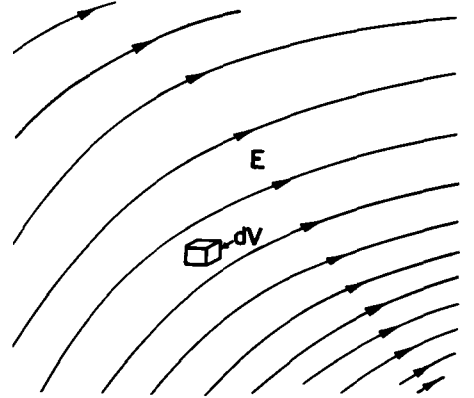


Fig. 8-8. Each volume element $dV = dx dy dz$ in an electric field contains the energy $(\epsilon_0/2)E^2 dV$.

charge in the distribution is zero.) Since the surface area of the large sphere increases as R^2 , we see that the surface integral falls off as $(1/R)(1/R^2)R^2 = (1/R)$ as the radius of the sphere increases. So if we include all space in our integration ($R \rightarrow \infty$), the surface integral goes to zero and we have that

$$U = \frac{\epsilon_0}{2} \int_{\text{all space}} (\nabla\phi) \cdot (\nabla\phi) dV = \frac{\epsilon_0}{2} \int_{\text{all space}} \mathbf{E} \cdot \mathbf{E} dV. \quad (8.35)$$

We see that it is possible for us to represent the energy of any charge distribution as being the integral over an energy density located in the field.

8-6 The energy of a point charge

Our new relation, Eq. (8.35), says that even a single point charge q will have some electrostatic energy. In this case, the electric field is given by

$$E = \frac{q}{4\pi\epsilon_0 r^2}.$$

So the energy density at the distance r from the charge is

$$\frac{\epsilon_0 E^2}{2} = \frac{q^2}{32\pi^2\epsilon_0 r^4}.$$

We can take for an element of volume a spherical shell of thickness dr and area $4\pi r^2$. The total energy is

$$U = \int_{r=0}^{\infty} \frac{q^2}{8\pi\epsilon_0 r^2} dr = -\frac{q^2}{8\pi\epsilon_0} \frac{1}{r} \Big|_{r=0}^{r=\infty}. \quad (8.36)$$

Now the limit at $r = \infty$ gives no difficulty. But for a point charge we are supposed to integrate down to $r = 0$, which gives an infinite integral. Equation (8.35) says that there is an infinite amount of energy in the field of a point charge, although we began with the idea that there was energy only *between* point charges. In our original energy formula for a collection of point charges (Eq. 8.3), we did not include any interaction energy of a charge with itself. What has happened is that when we went over to a continuous distribution of charge in Eq. (8.27), we counted the energy of interaction of every *infinitesimal* charge with all other infinitesimal charges. The same account is included in Eq. (8.35), so when we apply it to a *finite* point charge, we are including the energy it would take to assemble that charge from infinitesimal parts. You will notice, in fact, that we would also get the result in Eq. (8.36) if we used our expression (8.11) for the energy of a charged sphere and let the radius tend toward zero.

We must conclude that the idea of locating the energy in the field is inconsistent with the assumption of the existence of point charges. One way out of the difficulty would be to say that elementary charges, such as an electron, are not points but are really small distributions of charge. Alternatively, we could say that there is something wrong in our theory of electricity at very small distances, or with the idea of the local conservation of energy. There are difficulties with either point of view. These difficulties have never been overcome; they exist to this day. Sometime later, when we have discussed some additional ideas, such as the momentum in an electromagnetic field, we will give a more complete account of these fundamental difficulties in our understanding of nature.

Electricity in the Atmosphere

9-1 The electric potential gradient of the atmosphere

On an ordinary day over flat desert country, or over the sea, as one goes upward from the surface of the ground the electric potential increases by about 100 volts per meter. Thus there is a vertical electric field E of 100 volts/m in the air. The sign of the field corresponds to a negative charge on the earth's surface. This means that outdoors the potential at the height of your nose is 200 volts higher than the potential at your feet! You might ask: "Why don't we just stick a pair of electrodes out in the air one meter apart and use the 100 volts to power our electric lights?" Or you might wonder: "If there is *really* a potential difference of 200 volts between my nose and my feet, why is it I don't get a shock when I go out into the street?"

We will answer the second question first. Your body is a relatively good conductor. If you are in contact with the ground, you and the ground will tend to make one equipotential surface. Ordinarily, the equipotentials are parallel to the surface, as shown in Fig. 9-1(a), but when you are there, the equipotentials are distorted, and the field looks somewhat as shown in Fig. 9-1(b). So you still have very nearly zero potential difference between your head and your feet. There are charges that come from the earth to your head, changing the field. Some of them may be discharged by ions collected from the air, but the current of these is very small because air is a poor conductor.

9-1 The electric potential gradient of the atmosphere

9-2 Electric currents in the atmosphere

9-3 Origin of the atmospheric currents

9-4 Thunderstorms

9-5 The mechanism of charge separation

9-6 Lightning

Reference: Chalmers, J. Alan, *Atmospheric Electricity*, Pergamon Press, London (1957).

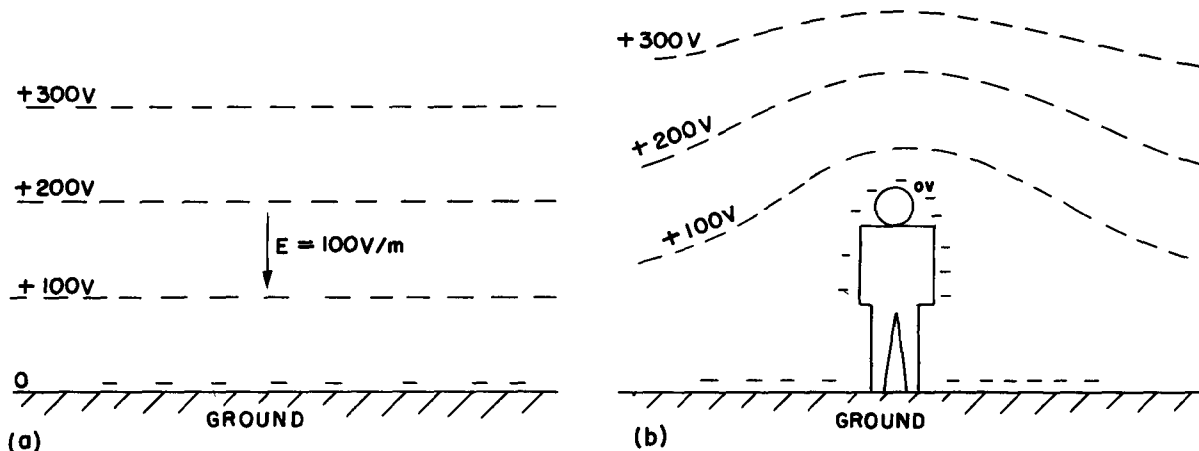


Fig. 9-1. (a) The potential distribution above the earth. (b) The potential distribution near a man in an open flat place.

How can we measure such a field if the field is changed by putting something there? There are several ways. One way is to place an insulated conductor at some distance above the ground and leave it there until it is at the same potential as the air. If we leave it long enough, the very small conductivity in the air will let the charges leak off (or onto) the conductor until it comes to the potential at its level. Then we can bring it back to the ground, and measure the shift of its potential as we do so. A faster way is to let the conductor be a bucket of water with a small leak. As the water drops out, it carries away any excess charges and the bucket will approach the same potential as the air. (The charges, as you know, reside on the surface, and as the drops come off "pieces of surface" break off.) We can measure the potential of the bucket with an electrometer.

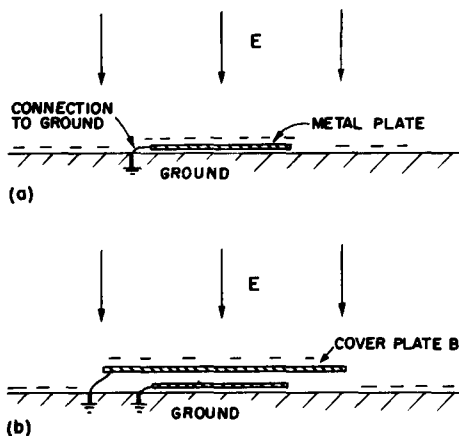


Fig. 9-2. (a) A grounded metal plate will have the same surface charge as the earth. (b) If the plate is covered with a grounded conductor it will have no surface charge.

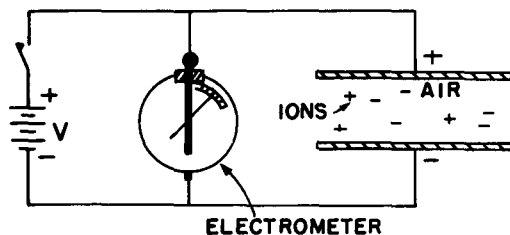


Fig. 9-3. Measuring the conductivity of air due to the motion of ions.

There is another way to directly measure the potential *gradient*. Since there is an electric field, there is a surface charge on the earth ($\sigma = \epsilon_0 E$). If we place a flat metal plate at the earth's surface and ground it, negative charges appear on it (Fig. 9-2a). If this plate is now covered by another grounded conducting cover *B*, the charges will appear on the cover, and there will be no charges on the original plate *A*. If we measure the charge that flows from plate *A* to the ground (by, say, a galvanometer in the grounding wire) as we cover it, we can find the surface charge density that was there, and therefore also find the electric field.

Having suggested how we can measure the electric field in the atmosphere, we now continue our description of it. Measurements show, first of all, that the field continues to exist, but gets weaker, as one goes up to high altitudes. By about 50 kilometers, the field is very small, so most of the potential change (the integral of E) is at lower altitudes. The total potential difference from the surface of the earth to the top of the atmosphere is about 400,000 volts.

9-2 Electric currents in the atmosphere

Another thing that can be measured, in addition to the potential gradient, is the current in the atmosphere. The current density is small—about 10 micromicro-amperes crosses each square meter parallel to the earth. The air is evidently not a perfect insulator, and because of this conductivity, a small current—caused by the electric field we have just been describing—passes from the sky down to the earth.

Why does the atmosphere have conductivity? Here and there among the air molecules there is an ion—a molecule of oxygen, say, which has acquired an extra electron, or perhaps lost one. These ions do not stay as single molecules; because of their electric field they usually accumulate a few other molecules around them. Each ion then becomes a little lump which, along with other lumps, drifts in the field—moving slowly upward or downward—making the observed current. Where do the *ions* come from? It was first guessed that the ions were produced by the radioactivity of the earth. (It was known that the radiation from radioactive materials would make air conducting by ionizing the air molecules.) Particles like β -rays coming out of the atomic nuclei are moving so fast that they tear electrons from the atoms, leaving ions behind. This would imply, of course, that if we were to go to higher altitudes, we should find less ionization, because the radioactivity is all in the dirt on the ground—in the traces of radium, uranium, potassium, etc.

To test this theory, some physicists carried an experiment up in balloons to measure the ionization of the air (Hess, in 1912) and discovered that the opposite was true—the ionization per unit volume *increased* with altitude! (The apparatus was like that of Fig. 9-3. The two plates were charged periodically to the potential V . Due to the conductivity of the air, the plates slowly discharged; the rate of discharge was measured with the electrometer.) This was a most mysterious result—the most dramatic finding in the entire history of atmospheric electricity. It was so dramatic, in fact, that it required a branching off of an entirely new subject—cosmic rays. Atmospheric electricity itself remained less dramatic. Ionization was evidently being produced by something from outside the earth; the investigation of this source led to the discovery of the cosmic rays. We will not discuss the subject of cosmic rays now, except to say that they maintain the supply of ions. Although the ions are being swept away all the time, new ones are being created by the cosmic-ray particles coming from the outside.

To be precise, we must say that besides the ions made of molecules, there are also other kinds of ions. Tiny pieces of dirt, like extremely fine bits of dust, float in the air and become charged. They are sometimes called “nuclei.” For example, when a wave breaks in the sea, little bits of spray are thrown into the air. When one of these drops evaporates, it leaves an infinitesimal crystal of NaCl floating in the air. These tiny crystals can then pick up charges and become ions; they are called “large ions.”

The small ions—those formed by cosmic rays—are the most mobile. Because they are so small, they move rapidly through the air—with a speed of about 1

cm/sec in a field of 100 volts/meter, or 1 volt/cm. The much bigger and heavier ions move much more slowly. It turns out that if there are many "nuclei," they will pick up the charges from the small ions. Then, since the "large ions" move so slowly in a field, the total conductivity is reduced. The conductivity of air, therefore, is quite variable, since it is very sensitive to the amount of "dirt" there is in it. There is much more of such dirt over land—where the winds can blow up dust or where man throws all kinds of pollution into the air—than there is over water. It is not surprising that from day to day, from moment to moment, from place to place, the conductivity near the earth's surface varies enormously. The voltage gradient observed at any particular place on the earth's surface also varies greatly because roughly the same current flows down from high altitudes in different places, and the varying conductivity near the earth results in a varying voltage gradient.

The conductivity of the air due to the drifting of ions also increases rapidly with altitude—for two reasons. First of all, the ionization from cosmic rays increases with altitude. Secondly, as the density of air goes down, the mean free path of the ions increases, so that they can travel farther in the electric field before they have a collision—resulting in a rapid increase of conductivity as one goes up.

Although the electric current-density in the air is only a few micromicroamperes per square meter, there are very many square meters on the earth's surface. The total electric current reaching the earth's surface at any time is very nearly constant at 1800 amperes. This current, of course, is "positive"—it carries plus charges to the earth. So we have a voltage supply of 400,000 volts with a current of 1800 amperes—a power of 700 megawatts!

With such a large current coming down, the negative charge on the earth should soon be discharged. In fact, it should take only about half an hour to discharge the entire earth. But the atmospheric electric field has already lasted more than a half-hour since its discovery. How is it maintained? What maintains the voltage? And between what and the earth? There are many questions.

The earth is negative, and the potential in the air is positive. If you go high enough, the conductivity is so great that horizontally there is no more chance for voltage variations. The air, for the scale of times that we are talking about, becomes effectively a conductor. This occurs at a height in the neighborhood of 50 kilometers. This is not as high as what is called the "ionosphere," in which there are very large numbers of ions produced by photoelectricity from the sun. Nevertheless, for our discussions of atmospheric electricity, the air becomes sufficiently conductive at about 50 kilometers that we can imagine that there is practically a perfect conducting surface at this height, from which the currents come down. Our picture of the situation is shown in Fig. 9-4. The problem is: How is the positive charge maintained there? How is it pumped back? Because if it comes down to the earth, it has to be pumped back somehow. That was one of the greatest puzzles of atmospheric electricity for quite a while.

Each piece of information we can get should give a clue or, at least, tell you something about it. Here is an interesting phenomenon: If we measure the current (which is more stable than the potential gradient) over the sea, for instance, or in careful conditions, and average very carefully so that we get rid of the irregularities, we discover that there is still a daily variation. The average of many measurements over the oceans has a variation with time roughly as shown in Fig. 9-5. The current varies by about ± 15 percent, and it is largest at 7:00 P.M. in London. The strange part of the thing is that no matter *where* you measure the current—in the Atlantic Ocean, the Pacific Ocean, or the Arctic Ocean—it is at its peak value when the clocks in *London* say 7:00 P.M.! All over the world the current is at its maximum at 7:00 P.M. London time and it is at a minimum at 4:00 A.M. London time. In other words, it depends upon the absolute time on the earth, *not* upon the local time at the place of observation. In one respect this is not mysterious; it checks with our idea that there is a very high conductivity laterally at the top, because that makes it impossible for the voltage difference from the ground to the top to vary locally. Any potential variations should be worldwide, as indeed they are. What we now know, therefore, is that the voltage at the "top" surface is dropping and rising by 15 percent with the absolute time on the earth.

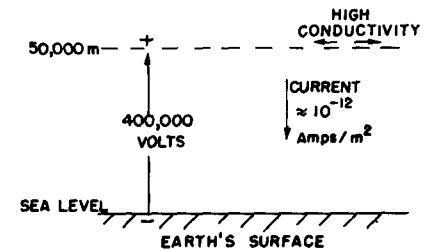


Fig. 9-4. Typical electrical conditions in a clear atmosphere.

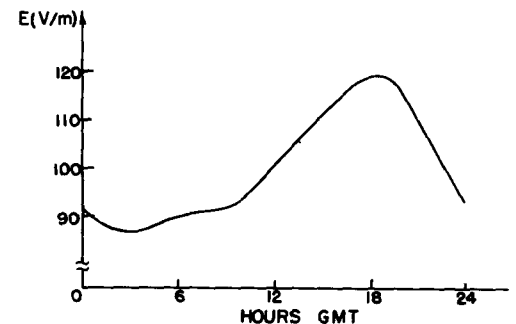


Fig. 9-5. The average daily variation of the atmospheric potential gradient on a clear day over the oceans; referred to Greenwich time.

9-3 Origin of the atmospheric currents

We must next talk about the source of the large negative currents which must be flowing from the “top” to the surface of the earth to keep charging it up negatively. Where are the batteries that do this? The “battery” is shown in Fig. 9-6. It is the thunderstorm and its lightning. It turns out that the bolts of lightning do not “discharge” the potential we have been talking about (as you might at first guess). Lightning storms carry *negative* charges to the earth. When a lightning bolt strikes, ten-to-one it brings down negative charges to the earth in large amounts. It is the thunderstorms throughout the world that are charging the earth with an average of 1800 amperes, which is then being discharged through regions of fair weather.

There are about 300 thunderstorms per day all over the earth, and we can think of them as batteries pumping the electricity to the upper layer and maintaining the voltage difference. Then take into account the geography of the earth—there are thunderstorms in the afternoon in Brazil, tropical thunderstorms in Africa, and so forth. People have made estimates of how much lightning is striking world-wide at any time, and perhaps needless to say, their estimates more or less agree with the voltage difference measurements: the total amount of thunderstorm activity is highest on the whole earth at about 7:00 P.M. in London. However, the thunderstorm estimates are very difficult to make and were made only *after* it was known that the variation should have occurred. These things are very difficult because we don't have enough observations on the seas and over all parts of the world to know the number of thunderstorms accurately. But those people who think they “do it right” obtain the result that there is a peak in the activity at 7:00 P.M. Greenwich Mean Time.



Fig. 9-6. The mechanism that generates the atmospheric electric field. [Photo by William L. Widmayer.]

In order to understand how these batteries work, we will look at a thunderstorm in detail. What is going on inside a thunderstorm? We will describe this insofar as it is known. As we get into this marvelous phenomenon of real nature—instead of the idealized spheres of perfect conductors inside of other spheres that we can solve so neatly—we discover that we don't know very much. Yet it is really quite exciting. Anyone who has been in a thunderstorm has enjoyed it, or has been frightened, or at least has had some emotion. And in those places in nature where we get an emotion, we find that there is generally a corresponding complexity and mystery about it. It is not going to be possible to describe exactly how a thunderstorm works, because we do not yet know very much. But we will try to describe a little bit about what happens.

9-4 Thunderstorms

In the first place, an ordinary thunderstorm is made up of a number of "cells" fairly close together, but almost independent of each other. So it is best to analyze one cell at a time. By a "cell" we mean a region with a limit area in the horizontal direction in which all of the basic processes occur. Usually there are several cells side by side, and in each one about the same thing is happening, although perhaps with a different timing. Figure 9-7 indicates in an idealized fashion what such a cell looks like in the early stage of the thunderstorm. It turns out that in a certain place in the air, under certain conditions which we shall describe, there is a general rising of the air, with higher and higher velocities near the top. As the warm, moist air at the bottom rises, it cools and condenses. In the figure the little crosses indicate snow and the dots indicate rain, but because the updraft currents are great enough and the drops are small enough, the snow and rain do not come down at this stage. This is the beginning stage, and not the real thunderstorm yet—in the sense that we don't have anything happening at the ground. At the same time that the warm air rises, there is an entrainment of air from the sides—an important point which was neglected for many years. Thus it is not just the air from below which is rising, but also a certain amount of other air from the sides.

Why does the air rise like this? As you know, when you go up in altitude the air is colder. The *ground* is heated by the sun, and the re-radiation of heat to the sky comes from water vapor high in the atmosphere; so at high altitudes the air is cold—very cold—whereas lower down it is warm. You may say, "Then it's very simple. Warm air is lighter than cold; therefore the combination is mechanically unstable and the warm air rises." Of course, if the temperature is different at different heights, the air *is* unstable *thermodynamically*. Left to itself infinitely long, the air would all come to the same temperature. But it is not left to itself; the sun is always shining (during the day). So the problem is indeed not one of thermodynamic equilibrium, but of *mechanical* equilibrium. Suppose we plot—as in Fig. 9-8—the temperature of the air against height above the ground. In ordinary circumstances we would get a decrease along a curve like the one labeled (a); as the height goes up, the temperature goes down. How can the atmosphere be stable? Why doesn't the hot air below simply rise up into the cold air? The answer is this: if the air were to go up, its pressure would go down, and if we consider a particular parcel of air going up, it would be expanding adiabatically. (There would be no heat coming in or out because in the large dimensions considered here, there isn't time for much heat flow.) Thus the parcel of air would cool as it rises. Such an adiabatic process would give a temperature-height relationship like curve (b) in Fig. 9-8. Any air which rose from below would be *colder* than the environment it goes into. Thus there is no reason for the hot air below to rise; if it were to rise, it would cool to a lower temperature than the air already there, would be heavier than the air there, and would just want to come down again. On a good, bright day with very little humidity there is a certain rate at which the temperature in the atmosphere falls, and this rate is, in general, lower than the "maximum stable gradient," which is represented by curve (b). The air is in stable mechanical equilibrium.

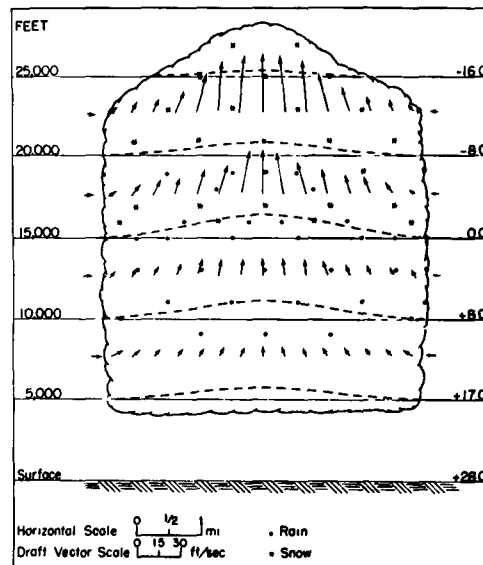


Fig. 9-7. A thunderstorm cell in the early stages of development. [From U.S. Department of Commerce Weather Bureau Report, June 1949.]

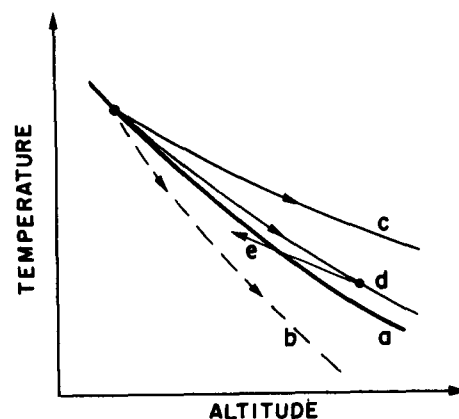


Fig. 9-8. Atmospheric temperature. (a) Static atmosphere; (b) adiabatic cooling of dry air; (c) adiabatic cooling of wet air; (d) wet air with some mixing of ambient air.

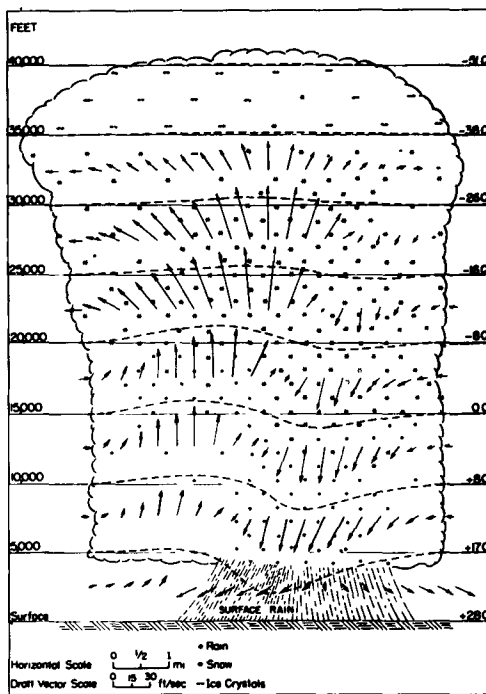


Fig. 9-9. A mature thunderstorm cell.
[From U.S. Department of Commerce
Weather Bureau Report, June 1949.]

On the other hand, if we think of a parcel of air that contains a lot of water vapor being carried up into the air, its adiabatic cooling curve will be different. As it expands and cools, the water vapor in it will condense, and the condensing water will liberate heat. Moist air, therefore, does not cool nearly as much as dry air does. So if air that is wetter than the average starts to rise, its temperature will follow a curve like (c) in Fig. 9-8. It will cool off somewhat, but will still be warmer than the surrounding air at the same level. If we have a region of warm moist air and something starts it rising, it will always find itself lighter and warmer than the air around it and will continue to rise until it gets to enormous heights. This is the machinery that makes the air in the thunderstorm cell rise.

For many years the thunderstorm cell was explained simply in this manner. But then measurements showed that the temperature of the cloud at different heights was not nearly as high as indicated by curve (c). The reason is that as the moist air "bubble" goes up, it entrains air from the environment and is cooled off by it. The temperature-versus-height curve looks more like curve (d), which is much closer to the original curve (a) than to curve (c).

After the convection just described gets under way, the cross section of a thunderstorm cell looks like Fig. 9-9. We have what is called a "mature" thunderstorm. There is a very rapid updraft which, in this stage, goes up to about 10,000 to 15,000 meters—sometimes even much higher. The thunderheads, with their condensation, climb way up out of the general cloud bank, carried by an updraft that is usually about 60 miles an hour. As the water vapor is carried up and condenses, it forms tiny drops which are rapidly cooled to temperatures below zero degrees. They should freeze, but do not freeze immediately—they are "supercooled." Water and other liquids will usually cool well below their freezing points before crystallizing if there are no "nuclei" present to start the crystallization process. Only if there is some small piece of material present, like a tiny crystal of NaCl, will the water drop freeze into a little piece of ice. Then the equilibrium is such that the water drops evaporate and the ice crystals grow. Thus at a certain point there is a rapid disappearance of the water and a rapid buildup of ice. Also, there may be direct collisions between the water drops and the ice—collisions in which the supercooled water becomes attached to the ice crystals, which causes it to suddenly crystallize. So at a certain point in the cloud expansion there is a rapid accumulation of large ice particles.

When the ice particles are heavy enough, they begin to fall through the rising air—they get too heavy to be supported any longer in the updraft. As they come down, they draw a little air with them and start a downdraft. And surprisingly enough, it is easy to see that once the downdraft is started, it will maintain itself. The air now drives itself down!

Notice that the curve (d) in Fig. 9-8 for the actual distribution of temperature in the cloud is not as steep as curve (c), which applies to wet air. So if we have wet air falling, its temperature will drop with the slope of curve (c) and will go *below* the temperature of the environment if it gets down far enough, as indicated by curve (e) in the figure. The moment it does that, it is denser than the environment and continues to fall rapidly. You say, "That is perpetual motion. First, you argue that the air should rise, and when you have it up there, you argue equally well that the air should fall." But it isn't perpetual motion. When the situation is unstable and the warm air should rise, then clearly something has to replace the warm air. It is equally true that cold air coming down would energetically replace the warm air, but you realize that what is coming down is *not* the original air. The early arguments, that had a particular cloud without entrainment going up and then coming down, had some kind of a puzzle. They needed the rain to maintain the downdraft—an argument which is hard to believe. As soon as you realize that there is a lot of original air mixed in with the rising air, the thermodynamic argument shows that there can be a descent of the cold air which was originally at some great height. This explains the picture of the active thunderstorm sketched in Fig. 9-9.

As the air comes down, rain begins to come out of the bottom of the thunderstorm. In addition, the relatively cold air spreads out when it arrives at the earth's surface. So just before the rain comes there is a certain little cold wind that gives

us a forewarning of the coming storm. In the storm itself there are rapid and irregular gusts of air, there is an enormous turbulence in the cloud, and so on. But basically we have an updraft, then a downdraft—in general, a very complicated process.

The moment at which precipitation starts is the same moment that the large downdraft begins and is the same moment, in fact, when the electrical phenomena arise. Before we describe lightning, however, we can finish the story by looking at what happens to the thunderstorm cell after about one-half an hour to an hour. The cell looks as shown in Fig. 9-10. The updraft stops because there is no longer enough warm air to maintain it. The downward precipitation continues for a while, the last little bits of water come out, and things get quieter and quieter—although there are small ice crystals left way up in the air. Because the winds at very great altitude are in different directions, the top of the cloud usually spreads into an anvil shape. The cell comes to the end of its life.

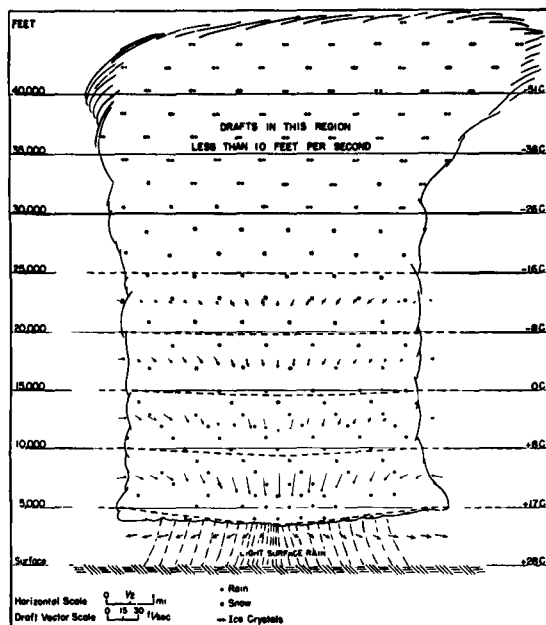


Fig. 9-10. The late phase of a thunderstorm cell. [From U.S. Department of Commerce Weather Bureau Report, June 1949.]

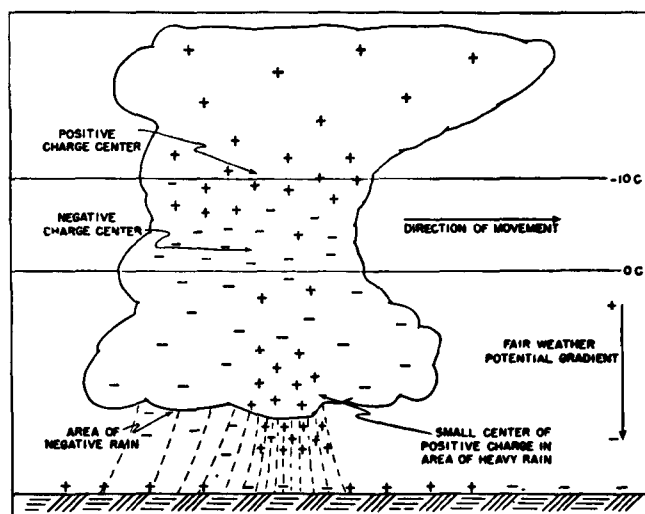


Fig. 9-11. The distribution of electrical charges in a mature thunderstorm cell. [From U.S. Department of Commerce Weather Bureau Report, June 1949.]

9-5 The mechanism of charge separation

We want now to discuss the most important aspect for our purposes—the development of the electrical charges. Experiments of various kinds—including flying airplanes through thunderstorms (the pilots who do this are brave men!)—tell us that the charge distribution in a thunderstorm cell is something like that shown in Fig. 9-11. The top of the thunderstorm has a positive charge, and the bottom a negative one—except for a small local region of positive charge in the bottom of the cloud, which has caused everybody a lot of worry. No one seems to know why it is there, how important it is—whether it is a secondary effect of the positive rain coming down, or whether it is an essential part of the machinery. Things would be much simpler if it weren't there. Anyway, the predominantly negative charge at the bottom and the positive charge at the top have the correct sign for the battery needed to drive the earth negative. The positive charges are 6 or 7 kilometers up in the air, where the temperature is about -20°C , whereas the negative charges are 3 or 4 kilometers high, where the temperature is between zero and -10°C .

The charge at the bottom of the cloud is large enough to produce potential differences of 20, or 30, or even 100 million volts between the cloud and the earth—much bigger than the 0.4 million volts from the “sky” to the ground in a clear

atmosphere. These large voltages break down the air and create giant arc discharges. When the breakdown occurs the negative charges at the bottom of the thunderstorm are carried down to the earth in the lightning strokes.

Now we will describe in some detail the character of the lightning. First of all, there are large voltage differences around, so that the air breaks down. There are lightning strokes between one piece of a cloud and another piece of a cloud, or between one cloud and another cloud, or between a cloud and the earth. In each of the independent discharge flashes—the kind of lightning strokes you see—there are approximately 20 or 30 coulombs of charge brought down. One question is: How long does it take for the cloud to regenerate the 20 or 30 coulombs which are taken away by the lightning bolt? This can be seen by measuring, far from a cloud, the electric field produced by the cloud's dipole moment. In such measurements you see a sudden decrease in the field when the lightning strikes, and then an exponential return to the previous value with a time constant which is slightly different for different cases but which is in the neighborhood of 5 seconds. It takes a thunderstorm only 5 seconds after each lightning stroke to build its charge up again. That doesn't necessarily mean that another stroke is going to occur in exactly 5 seconds every time, because, of course, the geometry is changed, and so on. The strokes occur more or less irregularly, but the important point is that it takes about 5 seconds to recreate the original condition. Thus there are approximately 4 amperes of current in the generating machine of the thunderstorm. This means that any model made to explain how this storm generates its electricity must be one with plenty of juice—it must be a big, rapidly operating device.

Before we go further we shall consider something which is almost certainly completely irrelevant, but nevertheless interesting, because it does show the effect of an electric field on water drops. We say that it may be irrelevant because it relates to an experiment one can do in the laboratory with a stream of water to show the rather strong effects of the electric field on drops of water. In a thunderstorm there is no stream of water; there is a cloud of condensing ice and drops of water. So the question of the mechanisms at work in a thunderstorm is probably not at all related to what you can see in the simple experiment we will describe. If you take a small nozzle connected to a water faucet and direct it upward at a steep angle, as in Fig. 9-12, the water will come out in a fine stream that eventually breaks up into a spray of fine drops. If you now put an electric field across the stream at the nozzle (by bringing up a charged rod, for example), the form of the stream will change. With a weak electric field you will find that the stream breaks up into a smaller number of large-sized drops. But if you apply a stronger field, the stream breaks up into many, many fine drops—smaller than before.* With a weak electric field there is a tendency to inhibit the breakup of the stream into drops. With a stronger field, however, there is an increase in the tendency to separate into drops.

The explanation of these effects is probably the following. If we have the stream of water coming out of the nozzle and we put a small electric field across it one side of the water gets slightly positive and the other side gets slightly negative. Then, when the stream breaks, the drops on one side may be positive, and those on the other side may be negative. They will attract each other and will have a tendency to stick together more than they would have before—the stream doesn't break up as much. On the other hand, if the field is stronger, the charge in each one of the drops gets much larger, and there is a tendency for the charge *itself* to help break up the drops through their own repulsion. Each drop will break into many smaller ones, each carrying a charge, so that they are all repelled, and spread out so rapidly. So as we increase the field, the stream becomes more finely separated. The only point we wish to make is that in certain circumstances electric fields can have considerable influence on the drops. The exact machinery by which something happens in a thunderstorm is not at all known, and is not at all necessarily related to what we have just described. We have included it just so that

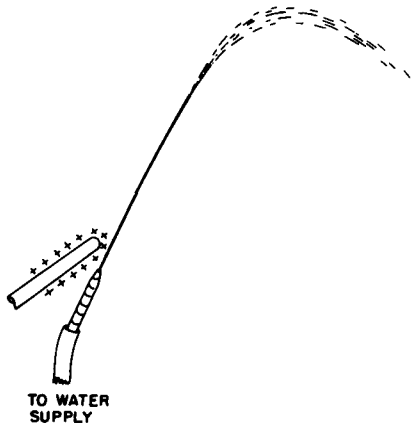


Fig. 9-12. A jet of water with an electric field near the nozzle.

* A handy way to observe the sizes of the drops is to let the stream fall on a large thin metal plate. The larger drops make a louder noise.

you will appreciate the complexities that could come into play. In fact, nobody has a theory applicable to clouds based on that idea.

We would like to describe two theories which have been invented to account for the separation of the charges in a thunderstorm. All the theories involve the idea that there should be some charge on the precipitation particles and a different charge in the air. Then by the movement of the precipitation particles—the water or the ice—through the air there is a separation of electric charge. The only question is: How does the charging of the drops begin? One of the older theories is called the “breaking-drop” theory. Somebody discovered that if you have a drop of water that breaks into two pieces in a windstream, there is positive charge on the water and negative charge in the air. This breaking-drop theory has several disadvantages, among which the most serious is that the *sign* is wrong. Second, in the large number of temperate-zone thunderstorms which do exhibit lightning, the precipitation effects at high altitudes are in ice, *not* in water.

From what we have just said, we note that if we could imagine some way for the charge to be different at the top and bottom of a drop and if we could also see some reason why drops in a high-speed airstream would break up into unequal pieces—a large one in the front and a smaller one in the back because of the motion through the air or something—we would have a theory. (Different from any known theory!) Then the small drops would not fall through the air as fast as the big ones, because of the air resistance, and we would get a charge separation. You see, it is possible to concoct all kinds of possibilities.

One of the more ingenious theories, which is more satisfactory in many respects than the breaking-drop theory, is due to C. T. R. Wilson. We will describe it, as Wilson did, with reference to water drops, although the same phenomenon would also work with ice. Suppose we have a water drop that is falling in the electric field of about 100 volts per meter toward the negatively charged earth. The drop will have an induced dipole moment—with the bottom of the drop positive and the top of the drop negative, as drawn in Fig. 9-13. Now there are in the air the “nuclei” that we mentioned earlier—the large slow-moving ions. (The fast ions do not have an important effect here.) Suppose that as a drop comes down, it approaches a large ion. If the ion is positive, it is repelled by the positive bottom of the drop and is pushed away. So it does not become attached to the drop. If the ion were to approach from the top, however, it might attach to the negative, top side. But since the drop is falling through the air, there is an air drift relative to it, going upwards, which carries the ions away if their motion through the air is slow enough. Thus the positive ions cannot attach at the top either. This would apply, you see, only to the large, slow-moving ions. The positive ions of this type will not attach themselves either to the front or the back of a falling drop. On the other hand, as the large, slow, *negative* ions are approached by a drop, they will be attracted and will be caught. The drop will acquire negative charge—the sign of the charge having been determined by the original potential difference on the entire earth—and we get the right sign. Negative charge will be brought down to the bottom part of the cloud by the drops, and the positively charged ions which are left behind will be blown to the top of the cloud by the various updraft currents. The theory looks pretty good, and it at least gives the right sign. Also it doesn't depend on having liquid drops. We will see, when we learn about polarization in a dielectric, that pieces of ice will do the same thing. They also will develop positive and negative charges on their extremities when they are in an electric field.

There are, however, some problems even with this theory. First of all, the total charge involved in a thunderstorm is very high. After a short time, the supply of large ions would get used up. So Wilson and others have had to propose that there are additional sources of the large ions. Once the charge separation starts, very large electric fields are developed, and in these large fields there may be places where the air will become ionized. If there is a highly charged point, or any small object like a drop, it may concentrate the field enough to make a “brush discharge.” When there is a strong enough electric field—let us say it is positive—electrons will fall into the field and will pick up a lot of speed between collisions. Their speed will be such that in hitting another atom they will tear electrons off at that

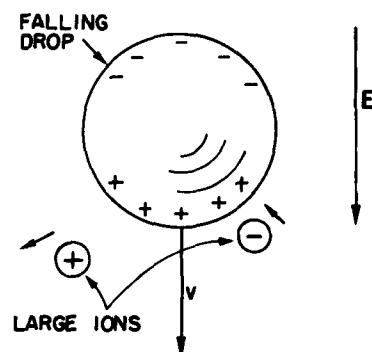


Fig. 9-13. C. T. R. Wilson's theory of charge separation in a thundercloud.

atom, leaving positive charges behind. These new electrons also pick up speed and collide with more electrons. So a kind of chain reaction or avalanche occurs, and there is a rapid accumulation of ions. The positive charges are left near their original positions, so the net effect is to distribute the positive charge on the point into a region around the point. Then, of course, there is no longer a strong field, and the process stops. This is the character of a brush discharge. It is possible that the fields may become strong enough in the cloud to produce a little bit of brush discharge; there may also be other mechanisms, once the thing is started, to produce a large amount of ionization. But nobody knows exactly how it works. So the fundamental origin of lightning is really not thoroughly understood. We know it comes from the thunderstorms. (And we know, of course, that thunder comes from the lightning—from the thermal energy released by the bolt.)

At least we can understand, in part, the origin of atmospheric electricity. Due to the air currents, ions, and water drops on ice particles in a thunderstorm, positive and negative charges are separated. The positive charges are carried upward to the top of the cloud (see Fig. 9-11), and the negative charges are dumped into the ground in lightning strokes. The positive charges leave the top of the cloud, enter the high-altitude layers of more highly conducting air, and spread throughout the earth. In regions of clear weather, the positive charges in this layer are slowly conducted to the earth by the ions in the air—ions formed by cosmic rays, by the sea, and by man's activities. The atmosphere is a busy electrical machine!

9-6 Lightning

The first evidence of what happens in a lightning stroke was obtained in photographs taken with a camera held by hand and moved back and forth with the shutter open—while pointed toward a place where lightning was expected. The first photographs obtained this way showed clearly that lightning strokes are usually multiple discharges along the same path. Later, the "Boys" camera, which has *two* lenses mounted 180° apart on a rapidly rotating disc, was developed. The image made by each lens moves across the film—the picture is spread out in time. If, for instance, the stroke repeats, there will be two images side by side. By comparing the images of the two lenses, it is possible to work out the details of the time sequence of the flashes. Figure 9-14 shows a photograph taken with a "Boys" camera.

We will now describe the lightning. Again, we don't understand exactly how it works. We will give a qualitative description of what it *looks* like, but we won't go into any details of *why* it does what it appears to do. We will describe only the ordinary case of the cloud with a negative bottom over flat country. Its potential is much more negative than the earth underneath, so negative electrons will be accelerated toward the earth. What happens is the following. It all starts with a thing called a "step leader," which is not as bright as the stroke of lightning. On the photographs one can see a little bright spot at the beginning that starts from the cloud and moves downward very rapidly—at a sixth of the speed of light! It goes only about 50 meters and stops. It pauses for about 50 microseconds, and then takes another step. It pauses again and then goes another step, and so on. It moves in a series of steps toward the ground, along a path like that shown in Fig. 9-15. In the leader there are negative charges from the cloud; the whole column is full of negative charge. Also, the air is becoming ionized by the rapidly moving charges that produce the leader, so the air becomes a conductor along the path traced out. The moment the leader touches the ground, we have a conducting "wire" that runs all the way up to the cloud and is full of negative charge. Now, at last, the negative charge of the cloud can simply escape and run out. The electrons at the bottom of the leader are the first ones to realize this; they dump out, leaving positive charge behind that attracts more negative charge from higher up in the leader, which in its turn pours out, etc. So finally all the negative charge in a part of the cloud runs out along the column in a rapid and energetic way. So the lightning stroke you *see* runs *upwards* from the ground, as indicated in Fig. 9-16. In fact, this main stroke—by far the brightest part—is called the *return*

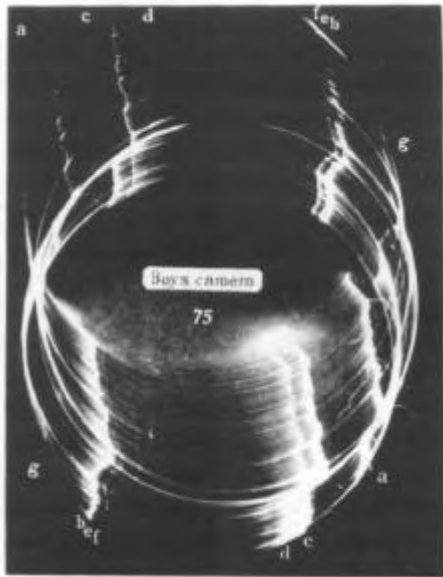


Fig. 9-14. Photograph of a lightning flash taken with a "Boys" camera. [From Schonland, Malan, and Collens, *Proc. Roy. Soc. London*, Vol. 152 (1935).]

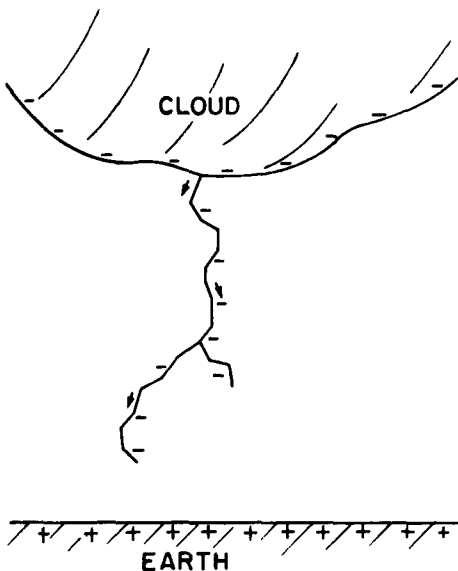


Fig. 9-15. The formation of the "step leader."

stroke. It is what produces the very bright light, and the heat, which by causing a rapid expansion of the air makes the thunder clap.

The current in a lightning stroke is about 10,000 amperes at its peak, and it carries down about 20 coulombs.

But we are still not finished. After a time of, perhaps, a few hundredths of a second, when the return stroke has disappeared, another leader comes down. But this time there are no pauses. It is called a “dark leader” this time, and it goes all the way down—from top to bottom in one swoop. It goes full steam on exactly the old track, because there is enough debris there to make it the easiest route. The new leader is again full of negative charge. The moment it touches the ground—zing!—there is a return stroke going straight up along the path. So you see the lightning strike again, and again, and again. Sometimes it strikes only once or twice, sometimes five or ten times—once as many as 42 times on the same track was seen—but always in rapid succession.

Sometimes things get even more complicated. For instance, after one of its pauses the leader may develop a branch by sending out *two* steps—both toward the ground but in somewhat different directions, as shown in Fig. 9-15. What happens then depends on whether one branch reaches the ground definitely before the other. If that does happen, the bright return stroke (of negative charge dumping into the ground) works its way *up* along the branch that touches the ground, and when it reaches and passes the branching point on its way up to the cloud, a bright stroke appears to go *down* the other branch. Why? Because negative charge is dumping out and that is what lights up the bolt. This charge begins to move at the top of the secondary branch, emptying successive, longer pieces of the branch, so the bright lightning bolt appears to work its way down that branch, at the same time as it works up toward the cloud. If, however, one of these extra leader branches happens to have reached the ground almost simultaneously with the original leader, it can sometimes happen that the *dark* leader of the second stroke will take the second branch. Then you will see the first main flash in one place and the second flash in another place. It is a variant of the original idea.

Also, our description is oversimplified for the region very near the ground. When the step leader gets to within a hundred meters or so from the ground, there is evidence that a discharge rises from the ground to meet it. Presumably, the field gets big enough for a brush-type discharge to occur. If, for instance, there is a sharp object, like a building with a point at the top, then as the leader comes down nearby the fields are so large that a discharge starts from the sharp point and reaches up to the leader. The lightning tends to strike such a point.

It has apparently been known for a long time that high objects are struck by lightning. There is a quotation of Artabanis, the advisor to Xerxes, giving his master advice on a contemplated attack on the Greeks—during Xerxes’ campaign to bring the entire known world under the control of the Persians. Artabanis said, “See how God with his lightning always smites the bigger animals and will not suffer them to wax insolent, while these of a lesser bulk chafe him not. How likewise his bolts fall ever on the highest houses and tallest trees.” And then he explains the reason: “So, plainly, doth he love to bring down everything that exalts itself.”

Do you think—now that you know a true account of lightning striking tall trees—that you have a greater wisdom in advising kings on military matters than did Artabanis 2300 years ago? Do not exalt yourself. You could only do it less poetically.

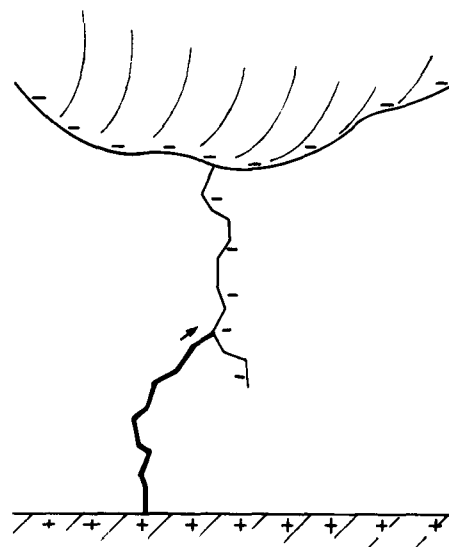


Fig. 9-16. The return lightning stroke runs back up the path made by the leader.

Dielectrics

10-1 The dielectric constant

Here we begin to discuss another of the peculiar properties of matter under the influence of the electric field. In an earlier chapter we considered the behavior of *conductors*, in which the charges move freely in response to an electric field to such points that there is no field left inside a conductor. Now we will discuss *insulators*, materials which do not conduct electricity. One might at first believe that there should be no effect whatsoever. However, using a simple electroscope and a parallel-plate capacitor, Faraday discovered that this was not so. His experiments showed that the capacitance of such a capacitor is *increased* when an insulator is put between the plates. If the insulator completely fills the space between the plates, the capacitance is increased by a factor κ which depends only on the nature of the insulating material. Insulating materials are also called *dielectrics*; the factor κ is then a property of the dielectric, and is called the *dielectric constant*. The dielectric constant of a vacuum is, of course, unity.

Our problem now is to explain why there is any electrical effect if the insulators are indeed insulators and do not conduct electricity. We begin with the experimental fact that the capacitance is increased and try to reason out what might be going on. Consider a parallel-plate capacitor with some charges on the surfaces of the conductors, let us say negative charge on the top plate and positive charge on the bottom plate. Suppose that the spacing between the plates is d and the area of each plate is A . As we have proved earlier, the capacitance is

$$C = \frac{\epsilon_0 A}{d}, \quad (10.1)$$

and the charge and voltage on the capacitor are related by

$$Q = CV. \quad (10.2)$$

Now the experimental fact is that if we put a piece of insulating material like lucite or glass between the plates, we find that the capacitance is larger. That means, of course, that the voltage is lower for the same charge. But the voltage difference is the integral of the electric field across the capacitor; so we must conclude that inside the capacitor, the electric field is reduced even though the charges on the plates remain unchanged.

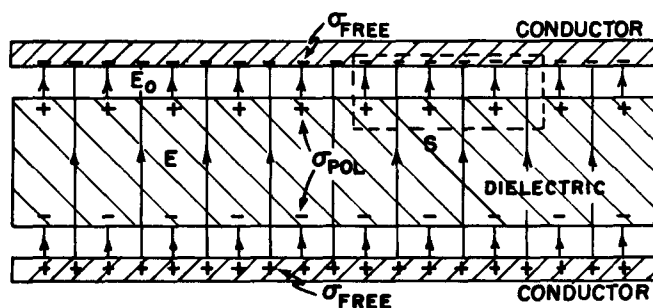


Fig. 10-1. A parallel-plate capacitor with a dielectric. The lines of E are shown.

Now how can that be? We have a law due to Gauss that tells us that the flux of the electric field is directly related to the enclosed charge. Consider the gaussian surface S shown by broken lines in Fig. 10-1. Since the electric field is reduced with the dielectric present, we conclude that the net charge inside the surface must

10-1 The dielectric constant

10-2 The polarization vector P

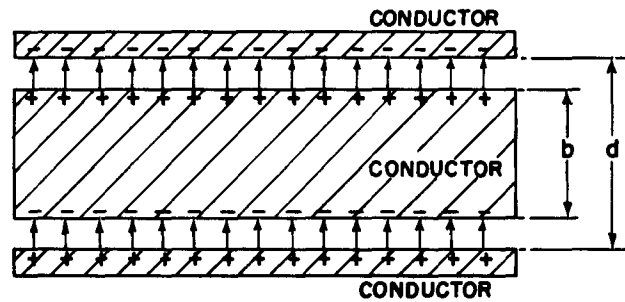
10-3 Polarization charges

10-4 The electrostatic equations with dielectrics

10-5 Fields and forces with dielectrics

be lower than it would be without the material. There is only one possible conclusion, and that is that there must be positive charges on the surface of the dielectric. Since the field is reduced but is not zero, we would expect this positive charge to be smaller than the negative charge on the conductor. So the phenomena can be explained if we could understand in some way that when a dielectric material is placed in an electric field there is positive charge induced on one surface and negative charge induced on the other.

Fig. 10-2. If we put a conducting plate in the gap of a parallel-plate condenser, the induced charges reduce the field in the conductor to zero.



We would expect that to happen for a conductor. For example, suppose that we had a capacitor with a plate spacing d , and we put between the plates a neutral conductor whose thickness is b , as in Fig. 10-2. The electric field induces a positive charge on the upper surface and a negative charge on the lower surface, so there is no field inside the conductor. The field in the rest of the space is the same as it was without the conductor, because it is the surface density of charge divided by ϵ_0 ; but the distance over which we have to integrate to get the voltage (the potential difference) is reduced. The voltage is

$$V = \frac{\sigma}{\epsilon_0} (d - b).$$

The resulting equation for the capacitance is like Eq. (10.1), with $(d - b)$ substituted for d :

$$C = \frac{\epsilon_0 A}{d[1 - (b/d)]}. \quad (10.3)$$

The capacitance is increased by a factor which depends upon (b/d) , the proportion of the volume which is occupied by the conductor.

This gives us an obvious model for what happens with dielectrics—that inside the material there are many little sheets of conducting material. The trouble with such a model is that it has a specific axis, the normal to the sheets, whereas most dielectrics have no such axis. However, this difficulty can be eliminated if we assume that all insulating materials contain small conducting spheres separated from each other by insulation, as shown in Fig. 10-3. The phenomenon of the dielectric constant is explained by the effect of the charges which would be induced on each sphere. This is one of the earliest physical models of dielectrics used to explain the phenomenon that Faraday observed. More specifically, it was assumed that each of the atoms of a material was a perfect conductor, but insulated from the others. The dielectric constant κ would depend on the proportion of space which was occupied by the conducting spheres. This is not, however, the model that is used today.

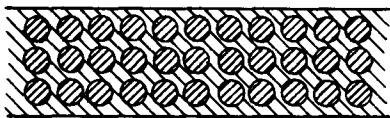


Fig. 10-3. A model of a dielectric: small conducting spheres embedded in an idealized insulator.

10-2 The polarization vector P

If we follow the above analysis further, we discover that the idea of regions of perfect conductivity and insulation is not essential. Each of the small spheres acts like a dipole, the moment of which is induced by the external field. The only thing that is essential to the understanding of dielectrics is that there are many little dipoles induced in the material. Whether the dipoles are induced because there are tiny conducting spheres or for any other reason is irrelevant.

Why should a field induce a dipole moment in an atom if the atom is not a conducting sphere? This subject will be discussed in much greater detail in the next chapter, which will be about the inner workings of dielectric materials. However, we give here one example to illustrate a possible mechanism. An atom has a positive charge on the nucleus, which is surrounded by negative electrons. In an electric field, the nucleus will be attracted in one direction and the electrons in the other. The orbits or wave patterns of the electrons (or whatever picture is used in quantum mechanics) will be distorted to some extent, as shown in Fig. 10-4; the center of gravity of the negative charge will be displaced and will no longer coincide with the positive charge of the nucleus. We have already discussed such distributions of charge. If we look from a distance, such a neutral configuration is equivalent, to a first approximation, to a little dipole.

It seems reasonable that if the field is not too enormous, the amount of induced dipole moment will be proportional to the field. That is, a small field will displace the charges a little bit and a larger field will displace them further—and in proportion to the field—unless the displacement gets too large. For the remainder of this chapter, it will be supposed that the dipole moment is exactly proportional to the field.

We will now assume that in each atom there are charges q separated by a distance δ , so that $q\delta$ is the dipole moment per atom. (We use δ because we are already using d for the plate separation.) If there are N atoms per unit volume, there will be a *dipole moment per unit volume* equal to $Nq\delta$. This dipole moment per unit volume will be represented by a vector, \mathbf{P} . Needless to say, it is in the direction of the individual dipole moments, i.e., in the direction of the charge separation δ :

$$\mathbf{P} = Nq\delta. \quad (10.4)$$

In general, \mathbf{P} will vary from place to place in the dielectric. However, at any point in the material, \mathbf{P} is proportional to the electric field \mathbf{E} . The constant of proportionality, which depends on the ease with which the electron are displaced, will depend on the kinds of atoms in the material.

What actually determines how this constant of proportionality behaves, how accurately it is constant for very large fields, and what is going on inside different materials, we will discuss at a later time. For the present, we will simply suppose that there exists a mechanism by which a dipole moment is induced which is proportional to the electric field.

10-3 Polarization charges

Now let us see what this model gives for the theory of a condenser with a dielectric. First consider a sheet of material in which there is a certain dipole moment per unit volume. Will there be on the average any charge density produced by this? Not if \mathbf{P} is uniform. If the positive and negative charges being displaced relative to each other have the same average density, the fact that they are displaced does not produce any net charge inside the volume. On the other hand, if \mathbf{P} were larger at one place and smaller at another, that would mean that more charge would be moved into some region than away from it; we would then expect to get a volume density of charge. For the parallel-plate condenser, we suppose that \mathbf{P} is uniform, so we need to look only at what happens at the surfaces. At one surface the negative charges, the electrons, have effectively moved out a distance δ ; at the other surface they have moved in, leaving some positive charge effectively out a distance δ . As shown in Fig. 10-5, we will have a surface density of charge, which will be called the *surface polarization charge*.

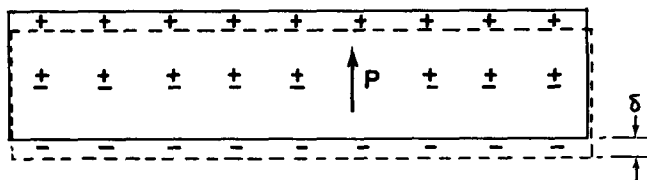


Fig. 10-5. A dielectric slab in a uniform field. The positive charges displaced the distance δ with respect to the negatives.

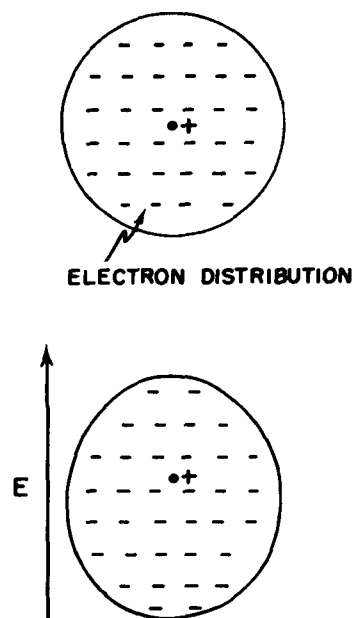


Fig. 10-4. An atom in an electric field has its distribution of electrons displaced with respect to the nucleus.

This charge can be calculated as follows. If A is the area of the plate, the number of electrons that appear at the surface is the product of A and N , the number per unit volume, and the displacement δ , which we assume here is perpendicular to the surface. The total charge is obtained by multiplying by the electronic charge q_e . To get the surface density of the polarization charge induced on the surface, we divide by A . The magnitude of the surface charge density is

$$\sigma_{\text{pol}} = Nq_e \delta.$$

But this is just equal to the magnitude P of the polarization vector \mathbf{P} , Eq. (10.4):

$$\sigma_{\text{pol}} = P. \quad (10.5)$$

The surface density of charge is equal to the polarization inside the material. The surface charge is, of course, positive on one surface and negative on the other.

Now let us assume that our slab is the dielectric of a parallel-plate capacitor. The *plates* of the capacitor also have a surface charge, which we will call σ_{free} , because they can move “freely” anywhere on the conductor. This is, of course, the charge that we put on when we charged the capacitor. It should be emphasized that σ_{pol} exists only because of σ_{free} . If σ_{free} is removed by discharging the capacitor, then σ_{pol} will disappear, not by going out on the discharging wire, but by moving back into the material—by the relaxation of the polarization inside the material.

We can now apply Gauss’ law to the gaussian surface S in Fig. 10–1. The electric field E in the dielectric is equal to the *total* surface charge density divided by ϵ_0 . It is clear that σ_{pol} and σ_{free} have opposite signs, so

$$E = \frac{\sigma_{\text{free}} - \sigma_{\text{pol}}}{\epsilon_0}. \quad (10.6)$$

Note that the field E_0 between the metal plate and the surface of the dielectric is higher than the field E ; it corresponds to σ_{free} alone. But here we are concerned with the field inside the dielectric which, if the dielectric nearly fills the gap, is the field over nearly the whole volume. Using Eq. (10.5), we can write

$$E = \frac{\sigma_{\text{free}} - P}{\epsilon_0}. \quad (10.7)$$

This equation doesn’t tell us what the electric field is unless we know what P is. Here, however, we are assuming that P depends on E —in fact, that it is proportional to E . This proportionality is usually written as

$$P = \chi \epsilon_0 E. \quad (10.8)$$

The constant χ (Greek “khi”) is called the *electric susceptibility* of the dielectric. Then Eq. (10.7) becomes

$$E = \frac{\sigma_{\text{free}}}{\epsilon_0} \frac{1}{(1 + \chi)}, \quad (10.9)$$

which gives us the factor $1/(1 + \chi)$ by which the field is reduced.

The voltage between the plates is the integral of the electric field. Since the field is uniform, the integral is just the product of E and the plate separation d . We have that

$$V = Ed = \frac{\sigma_{\text{free}} d}{\epsilon_0(1 + \chi)}.$$

The total charge on the capacitor is $\sigma_{\text{free}} A$, so that the capacitance defined by (10.2) becomes

$$C = \frac{\epsilon_0 A(1 + \chi)}{d} = \frac{\kappa \epsilon_0 A}{d}. \quad (10.10)$$

We have explained the observed facts. When a parallel-plate capacitor is filled with a dielectric, the capacitance is increased by the factor

$$\kappa = 1 + \chi, \quad (10.11)$$

which is a property of the material. Our explanation, of course, is not complete until we have explained—as we will do later—how the atomic polarization comes about.

Let's now consider something a little bit more complicated—the situation in which the polarization \mathbf{P} is not everywhere the same. As mentioned earlier, if the polarization is not constant, we would expect in general to find a charge density in the volume, because more charge might come into one side of a small volume element than leaves it on the other. How can we find out how much charge is gained or lost from a small volume?

First let's compute how much charge moves across any imaginary surface when the material is polarized. The amount of charge that goes across a surface is just P times the surface area if the polarization is *normal* to the surface. Of course, if the polarization is *tangential* to the surface, no charge moves across it.

Following the same arguments we have already used, it is easy to see that the charge moved across any surface element is proportional to the *component* of \mathbf{P} *perpendicular* to the surface. Compare Fig. 10-6 with Fig. 10-5. We see that Eq. (10.5) should, in the general case, be written

$$\sigma_{\text{pol}} = \mathbf{P} \cdot \mathbf{n}. \quad (10.12)$$

If we are thinking of an imagined surface element *inside* the dielectric, Eq. (10.12) gives the charge moved across the surface but doesn't result in a net surface charge, because there are equal and opposite contributions from the dielectric on the two sides of the surface.

The displacements of the charges can, however, result in a *volume* charge density. The total charge displaced *out* of any volume V by the polarization is the integral of the outward normal component of \mathbf{P} over the surface S that bounds the volume (see Fig. 10-7). An equal excess charge of the opposite sign is left behind. Denoting the net charge inside V by ΔQ_{pol} we write

$$\Delta Q_{\text{pol}} = - \int_S \mathbf{P} \cdot \mathbf{n} \, da. \quad (10.13)$$

We can attribute ΔQ_{pol} to a volume distribution of charge with the density ρ_{pol} , and so

$$\Delta Q_{\text{pol}} = \int_V \rho_{\text{pol}} \, dV. \quad (10.14)$$

Combining the two equations yields

$$\int_V \rho_{\text{pol}} \, dV = - \int_S \mathbf{P} \cdot \mathbf{n} \, da. \quad (10.15)$$

We have a kind of Gauss' theorem that relates the charge density from polarized materials to the polarization vector \mathbf{P} . We can see that it agrees with the result we got for the surface polarization charge or the dielectric in a parallel-plate capacitor. Using Eq. (10.15) with the gaussian surface of Fig. 10-1, the surface integral gives $P \Delta A$, and the charge inside is $\sigma_{\text{pol}} \Delta A$, so we get again that $\sigma = P$.

Just as we did for Gauss' law of electrostatics, we can convert Eq. (10.15) to a differential form—using Gauss' mathematical theorem:

$$\int_S \mathbf{P} \cdot \mathbf{n} \, da = \int_V \nabla \cdot \mathbf{P} \, dV.$$

We get

$$\rho_{\text{pol}} = - \nabla \cdot \mathbf{P}. \quad (10.16)$$

If there is a nonuniform polarization, its divergence gives the net density of charge appearing in the material. We emphasize that this is a perfectly *real* charge density; we call it “polarization charge” only to remind ourselves how it got there.

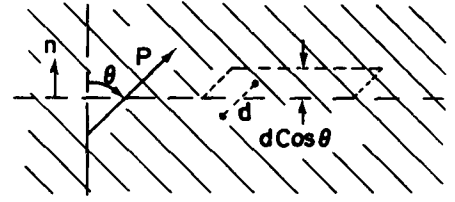


Fig. 10-6. The charge moved across an element of an imaginary surface in a dielectric is proportional to the component of \mathbf{P} normal to the surface.

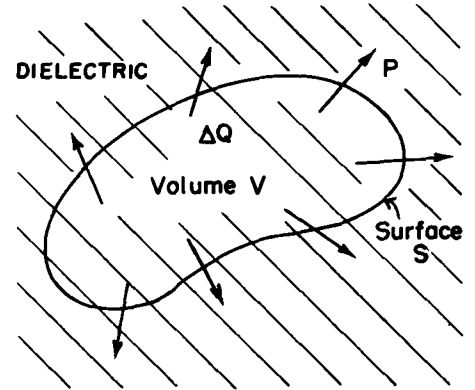


Fig. 10-7. A nonuniform polarization \mathbf{P} can result in a net charge in the body of a dielectric.

10-4 The electrostatic equations with dielectrics

Now let's combine the above result with our theory of electrostatics. The fundamental equation is

$$\nabla \cdot \mathbf{E} = \frac{\rho}{\epsilon_0}. \quad (10.17)$$

The ρ here is the density of *all* electric charges. Since it is not easy to keep track of the polarization charges, it is convenient to separate ρ into two parts. Again we call ρ_{pol} the charges due to nonuniform polarizations, and call ρ_{free} all the rest. Usually ρ_{free} is the charge we put on conductors, or at known places in space. Equation (10.17) then becomes

$$\nabla \cdot \mathbf{E} = \frac{\rho_{\text{free}} + \rho_{\text{pol}}}{\epsilon_0} = \frac{\rho_{\text{free}} - \nabla \cdot \mathbf{P}}{\epsilon_0},$$

or

$$\nabla \cdot \left(\mathbf{E} + \frac{\mathbf{P}}{\epsilon_0} \right) = \frac{\rho_{\text{free}}}{\epsilon_0}. \quad (10.18)$$

Of course, the equation for the curl of \mathbf{E} is unchanged:

$$\nabla \times \mathbf{E} = 0. \quad (10.19)$$

Taking \mathbf{P} from Eq. (10.8), we get the simpler equation

$$\nabla \cdot [(1 + \chi)\mathbf{E}] = \nabla \cdot (\kappa\mathbf{E}) = \frac{\rho_{\text{free}}}{\epsilon_0}. \quad (10.20)$$

These are the equations of electrostatics when there are dielectrics. They don't, of course, say anything new, but they are in a form which is more convenient for computation in cases where ρ_{free} is known and the polarization \mathbf{P} is proportional to \mathbf{E} .

Notice that we have not taken the dielectric "constant," κ , out of the divergence. That is because it may not be the same everywhere. If it has everywhere the same value, it can be factored out and the equations are just those of electrostatics with the charge density ρ_{free} divided by κ . In the form we have given, the equations apply to the general case where different dielectrics may be in different places in the field. Then the equations may be quite difficult to solve.

There is a matter of some historical importance which should be mentioned here. In the early days of electricity, the atomic mechanism of polarization was not known and the existence of ρ_{pol} was not appreciated. The charge ρ_{free} was considered to be the entire charge density. In order to write Maxwell's equations in a simple form, a new vector \mathbf{D} was defined to be equal to a linear combination of \mathbf{E} and \mathbf{P} :

$$\mathbf{D} = \epsilon_0 \mathbf{E} + \mathbf{P}. \quad (10.21)$$

As a result, Eqs. (10.18) and (10.19) were written in an apparently very simple form:

$$\nabla \cdot \mathbf{D} = \rho_{\text{free}}, \quad \nabla \times \mathbf{E} = 0. \quad (10.22)$$

Can one solve these? Only if a third equation is given for the relationship between \mathbf{D} and \mathbf{E} . When Eq. (10.8) holds, this relationship is

$$\mathbf{D} = \epsilon_0(1 + \chi)\mathbf{E} = \kappa\epsilon_0\mathbf{E}. \quad (10.23)$$

This equation was usually written

$$\mathbf{D} = \epsilon\mathbf{E}, \quad (10.24)$$

where ϵ is still another constant for describing the dielectric property of materials. It is called the "permittivity." (Now you see why we have ϵ_0 in our equations, it is the "permittivity of empty space.") Evidently,

$$\epsilon = \kappa\epsilon_0 = (1 + \chi)\epsilon_0. \quad (10.25)$$

Today we look upon these matters from another point of view, namely, that we have simpler equations in a vacuum, and if we exhibit in every case all the charges, whatever their origin, the equations are always correct. If we separate some of the charges away for convenience, or because we do not want to discuss what is going on in detail, then we can, if we wish, write our equations in any other form that may be convenient.

One more point should be emphasized. An equation like $\mathbf{D} = \epsilon \mathbf{E}$ is an attempt to describe a property of matter. But matter is extremely complicated, and such an equation is in fact not correct. For instance, if E gets too large, then \mathbf{D} is no longer proportional to \mathbf{E} . For some substances, the proportionality breaks down even with relatively small fields. Also, the "constant" of proportionality may depend on how fast E changes with time. Therefore this kind of equation is a kind of approximation, like Hooke's law. It cannot be a deep and fundamental equation. On the other hand, our fundamental equations for \mathbf{E} , (10.17) and (10.19), represent our deepest and most complete understanding of electrostatics.

10-5 Fields and forces with dielectrics

We will now prove some rather general theorems for electrostatics in situations where dielectrics are present. We have seen that the capacitance of a parallel-plate capacitor is increased by a definite factor if it is filled with a dielectric. We can show that this is true for a capacitor of *any* shape, provided the entire region in the neighborhood of the two conductors is filled with a uniform linear dielectric. Without the dielectric, the equations to be solved are

$$\nabla \cdot \mathbf{E}_0 = \frac{\rho_{\text{free}}}{\epsilon_0} \quad \text{and} \quad \nabla \times \mathbf{E}_0 = 0.$$

With the dielectric present, the first of these equations is modified; we have instead the equations

$$\nabla \cdot (\kappa \mathbf{E}) = \frac{\rho_{\text{free}}}{\epsilon_0} \quad \text{and} \quad \nabla \times \mathbf{E} = 0. \quad (10.26)$$

Now since we are taking κ to be everywhere the same, the last two equations can be written as

$$\nabla \cdot (\kappa \mathbf{E}) = \frac{\rho_{\text{free}}}{\epsilon_0} \quad \text{and} \quad \nabla \times (\kappa \mathbf{E}) = 0. \quad (10.27)$$

We therefore have the same equations for $\kappa \mathbf{E}$ as for \mathbf{E}_0 , so they have the solution $\kappa \mathbf{E} = \mathbf{E}_0$. In other words, the field is everywhere smaller, by the factor $1/\kappa$, than in the case without the dielectric. Since the voltage difference is a line integral of the field, the voltage is reduced by this same factor. Since the charge on the electrodes of the capacitor has been taken the same in both cases, Eq. (10.2) tells us that the capacitance, in the case of an everywhere uniform dielectric, is increased by the factor κ .

Let us now ask what the *force* would be between two charged conductors in a dielectric. We consider a liquid dielectric that is homogeneous everywhere. We have seen earlier that one way to obtain the force is to differentiate the energy with respect to the appropriate distance. If the conductors have equal and opposite charges, the energy $U = Q^2/2C$, where C is their capacitance. Using the principle of virtual work, any component is given by a differentiation; for example,

$$F_x = -\frac{\partial U}{\partial x} = -\frac{Q^2}{2} \frac{\partial}{\partial x} \left(\frac{1}{C} \right). \quad (10.28)$$

Since the dielectric increases the capacity by a factor κ , all forces will be *reduced* by this same factor.

One point should be emphasized. What we have said is true only if the dielectric is a liquid. Any motion of conductors that are embedded in solid dielectric changes the mechanical stress conditions of the dielectric and alters its electrical

properties, as well as causing some mechanical energy change in the dielectric. Moving the conductors in a liquid does not change the liquid. The liquid moves to a new place but its electrical characteristics are not changed.

Many older books on electricity start with the “fundamental” law that the force between two charges is

$$F = \frac{q_1 q_2}{4\pi\epsilon_0 \kappa r^2}, \quad (10.29)$$

a point of view which is thoroughly unsatisfactory. For one thing, it is not true in general; it is true only for a world filled with a liquid. Secondly, it depends on the fact that κ is a constant, which is only approximately true for most real materials. It is much better to start with Coulomb’s law for charges in a *vacuum*, which is always right (for stationary charges).

What does happen in a solid? This is a very difficult problem which has not been solved, because it is, in a sense, indeterminate. If you put charges inside a dielectric solid, there are many kinds of pressures and strains. You cannot deal with virtual work without including also the mechanical energy required to compress the solid, and it is a difficult matter, generally speaking, to make a unique distinction between the electrical forces and the mechanical forces due to the solid material itself. Fortunately, no one ever really needs to know the answer to the question proposed. He may sometimes want to know how much strain there is going to be in a solid, and that can be worked out. But it is much more complicated than the simple result we got for liquids.

A surprisingly complicated problem in the theory of dielectrics is the following: Why does a charged object pick up little pieces of dielectric? If you comb your hair on a dry day, the comb readily picks up small scraps of paper. If you thought casually about it, you probably assumed the comb had one charge on it and the paper had the opposite charge on it. But the paper is initially electrically neutral. It hasn’t any net charge, but it is attracted anyway. It is true that sometimes the paper will come up to the comb and then fly away, repelled immediately after it touches the comb. The reason is, of course, that when the paper touches the comb, it picks up some negative charges and then the like charges repel. But that doesn’t answer the original question. Why did the paper come toward the comb in the first place?

The answer has to do with the polarization of a dielectric when it is placed in an electric field. There are polarization charges of both signs, which are attracted and repelled by the comb. There is a net attraction, however, because the field nearer the comb is stronger than the field farther away—the comb is not an infinite sheet. Its charge is localized. A neutral piece of paper will not be attracted to either plate inside the parallel plates of a capacitor. The variation of the field is an essential part of the attraction mechanism.

As illustrated in Fig. 10–8, a dielectric is always drawn from a region of weak field toward a region of stronger field. In fact, one can prove that for small objects the force is proportional to the gradient of the *square* of the electric field. Why does it depend on the square of the field? Because the induced polarization charges are proportional to the fields, and for given charges the forces are proportional to the field. However, as we have just indicated, there will be a *net* force only if the square of the field is changing from point to point. So the force is proportional to the gradient of the square of the field. The constant of proportionality involves, among other things, the dielectric constant of the object, and it also depends upon the size and shape of the object.

There is a related problem in which the force on a dielectric can be worked out quite accurately. If we have a parallel-plate capacitor with a dielectric slab only partially inserted, as shown in Fig. 10–9, there will be a force driving the sheet in. A detailed examination of the force is quite complicated; it is related to nonuniformities in the field near the edges of the dielectric and the plates. However, if we do not look at the details, but merely use the principle of conservation of energy, we can easily calculate the force. We can find the force from the formula we de-

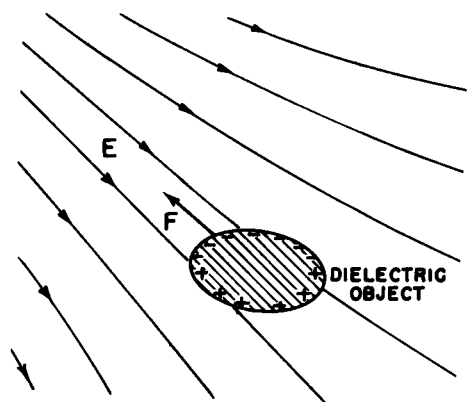


Fig. 10–8. A dielectric object in a nonuniform field feels a force toward regions of higher field strength.

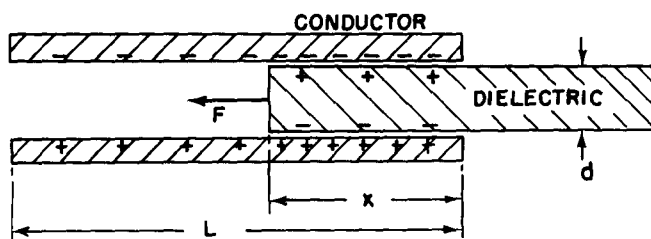


Fig. 10-9. The force on a dielectric sheet in a parallel-plate capacitor can be computed by applying the principle of energy conservation.

rived earlier. Equation (10.28) is equivalent to

$$F_x = -\frac{\partial U}{\partial x} = +\frac{V^2}{2} \frac{\partial C}{\partial x}. \quad (10.30)$$

We need only find out how the capacitance varies with the position of the dielectric slab.

Let's suppose that the total length of the plates is L , that the width of the plates is W , that the plate separation and dielectric thickness are d , and that the distance to which the dielectric has been inserted is x . The capacitance is the ratio of the total free charge on the plates to the voltage between the plates. We have seen above that for a given voltage V the surface charge density of free charge is $\kappa\epsilon_0 V/d$. So the total charge on the plates is

$$Q = \frac{\kappa\epsilon_0 V}{d} xW + \frac{\epsilon_0 V}{d} (L - x)W,$$

from which we get the capacitance:

$$C = \frac{\epsilon_0 W}{d} (\kappa x + L - x). \quad (10.31)$$

Using (10.30), we have

$$F_x = \frac{V^2}{2} \frac{\epsilon_0 W}{d} (\kappa - 1). \quad (10.32)$$

Now this equation is not particularly useful for anything unless you happen to need to know the force in such circumstances. We only wished to show that the theory of energy can often be used to avoid enormous complications in determining the forces on dielectric materials—as there would be in the present case.

Our discussion of the theory of dielectrics has dealt only with electrical phenomena, accepting the fact that the material has a polarization which is proportional to the electric field. Why there is such a proportionality is perhaps of greater interest to physics. Once we understand the origin of the dielectric constants from an atomic point of view, we can use electrical measurements of the dielectric constants in varying circumstances to obtain detailed information about atomic or molecular structure. This aspect will be treated in part in the next chapter.

Inside Dielectrics

11-1 Molecular dipoles

In this chapter we are going to discuss why it is that materials are dielectric. We said in the last chapter that we could understand the properties of electrical systems with dielectrics once we appreciated that when an electric field is applied to a dielectric it induces a dipole moment in the atoms. Specifically, if the electric field E induces an average dipole moment per unit volume P , then κ , the dielectric constant, is given by

$$\kappa - 1 = \frac{P}{\epsilon_0 E}. \quad (11.1)$$

We have already discussed how this equation is applied; now we have to discuss the mechanism by which polarization arises when there is an electric field inside a material. We begin with the simplest possible example—the polarization of gases. But even gases already have complications: there are two types. The molecules of some gases, like oxygen, which has a symmetric pair of atoms in each molecule, have no inherent dipole moment. But the molecules of others, like water vapor (which has a nonsymmetric arrangement of hydrogen and oxygen atoms) carry a permanent electric dipole moment. As we pointed out in Chapters 6 and 7, there is in the water vapor molecule an average plus charge on the hydrogen atoms and a negative charge on the oxygen. Since the center of gravity of the negative charge and the center of gravity of the positive charge do not coincide, the total charge distribution of the molecule has a dipole moment. Such a molecule is called a *polar* molecule. In oxygen, because of the symmetry of the molecule, the centers of gravity of the positive and negative charges are the same, so it is a *nonpolar* molecule. It does, however, become a dipole when placed in an electric field. The forms of the two types of molecules are sketched in Fig. 11-1.

11-2 Electronic polarization

We will first discuss the polarization of nonpolar molecules. We can start with the simplest case of a monatomic gas (for instance, helium). When an atom of such a gas is in an electric field, the electrons are pulled one way by the field while the nucleus is pulled the other way, as shown in Fig. 10-4. Although the atoms are very stiff with respect to the electrical forces we can apply experimentally, there is a slight net displacement of the centers of charge, and a dipole moment is induced. For small fields, the amount of displacement, and so also the dipole moment, is proportional to the electric field. The displacement of the electron distribution which produces this kind of induced dipole moment is called *electronic polarization*.

We have already discussed the influence of an electric field on an atom in Chapter 31 of Vol. I, when we were dealing with the theory of the index of refraction. If you think about it for a moment, you will see that what we must do now is exactly the same as we did then. But now we need worry only about fields that do not vary with time, while the index of refraction depended on time-varying fields.

In Chapter 31 of Vol. I we supposed that when an atom is placed in an oscillating electric field the center of charge of the electrons obeys the equation

$$m \frac{d^2 x}{dt^2} + m\omega_0^2 x = q_e E. \quad (11.2)$$

11-1 Molecular dipoles

11-2 Electronic polarization

11-3 Polar molecules; orientation polarization

11-4 Electric fields in cavities of a dielectric

11-5 The dielectric constant of liquids; the Clausius-Mossotti equation

11-6 Solid dielectrics

11-7 Ferroelectricity; BaTiO₃

Review: Chapter 31, Vol. I, *The Origin of the Refractive Index*
Chapter 40, Vol. I, *The Principles of Statistical Mechanics*

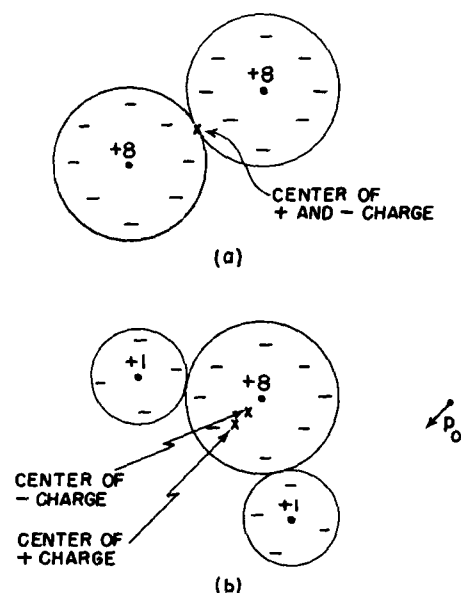


Fig. 11-1. (a) An oxygen molecule with zero dipole moment. (b) The water molecule has a permanent dipole moment p_0 .

The first term is the electron mass times its acceleration and the second is a restoring force, while the right-hand side is the force from the outside electric field. If the electric field varies with the frequency ω , Eq. (11.2) has the solution

$$x = \frac{q_e E}{m(\omega_0^2 - \omega^2)}, \quad (11.3)$$

which has a resonance at $\omega = \omega_0$. When we previously found this solution, we interpreted it as saying that ω_0 was the frequency at which light (in the optical region or in the ultraviolet, depending on the atom) was absorbed. For our purposes, however, we are interested only in the case of constant fields, i.e., for $\omega = 0$, so we can disregard the acceleration term in (11.2), and we find that the displacement is

$$x = \frac{q_e E}{m\omega_0^2}. \quad (11.4)$$

From this we see that the dipole moment p of a single atom is

$$p = q_e x = \frac{q_e^2 E}{m\omega_0^2}. \quad (11.5)$$

In this theory the dipole moment p is indeed proportional to the electric field.

People usually write

$$p = \alpha \epsilon_0 E. \quad (11.6)$$

(Again the ϵ_0 is put in for historical reasons.) The constant α is called the polarizability of the atom, and has the dimensions L^3 . It is a measure of how easy it is to induce a moment in an atom with an electric field. Comparing (11.5) and (11.6), our simple theory says that

$$\alpha = \frac{q_e^2}{\epsilon_0 m\omega_0^2} = \frac{4\pi e^2}{m\omega_0^2}. \quad (11.7)$$

If there are N atoms in a unit volume, the polarization P —the dipole moment per unit volume—is given by

$$P = Np = N\alpha \epsilon_0 E. \quad (11.8)$$

Putting (11.1) and (11.8) together, we get

$$\kappa - 1 = \frac{P}{\epsilon_0 E} = N\alpha \quad (11.9)$$

or, using (11.7),

$$\kappa - 1 = \frac{4\pi N e^2}{m\omega_0^2}. \quad (11.10)$$

From Eq. (11.9) we would predict that the dielectric constant κ of different gases should depend on the density of the gas and on the frequency ω_0 of its optical absorption.

Our formula is, of course, only a very rough approximation, because in Eq. (11.2) we have taken a model which ignores the complications of quantum mechanics. For example, we have assumed that an atom has only one resonant frequency, when it really has many. To calculate properly the polarizability α of atoms we must use the complete quantum-mechanical theory, but the classical ideas above give us a reasonable estimate.

Let's see if we can get the right order of magnitude for the dielectric constant of some substance. Suppose we try hydrogen. We have once estimated (Chapter 38, Vol. I) that the energy needed to ionize the hydrogen atom should be approximately

$$E \approx \frac{1}{2} \frac{me^4}{\hbar^2}. \quad (11.11)$$

For an estimate of the natural frequency ω_0 , we can set this energy equal to $\hbar\omega_0$ —the energy of an atomic oscillator whose natural frequency is ω_0 . We get

$$\omega_0 \approx \frac{1}{2} \frac{me^4}{\hbar^3}.$$

If we now use this value of ω_0 in Eq. (11.7), we find for the electronic polarizability

$$\alpha \approx 16\pi \left[\frac{\hbar^2}{me^2} \right]^3. \quad (11.12)$$

The quantity (\hbar^2/me^2) is the radius of the ground-state orbit of a Bohr atom (see Chapter 38, Vol. I) and equals 0.528 angstroms. In a gas at standard pressure and temperature (1 atmosphere, 0°C) there are 2.69×10^{19} atoms/cm³, so Eq. (11.9) gives us

$$\kappa = 1 + (2.69 \times 10^{19})16\pi(0.528 \times 10^{-8})^3 = 1.00020. \quad (11.13)$$

The dielectric constant for hydrogen gas is measured to be

$$\kappa_{\text{exp}} = 1.00026.$$

We see that our theory is about right. We should not expect any better, because the measurements were, of course, made with normal hydrogen gas, which has diatomic molecules, not single atoms. We should not be surprised if the polarization of the atoms in a molecule is not quite the same as that of the separate atoms. The molecular effect, however, is not really that large. An exact quantum-mechanical calculation of α for hydrogen atoms gives a result about 12% higher than (11.12) (the 16π is changed to 18π), and therefore predicts a dielectric constant somewhat closer to the observed one. In any case, it is clear that our model of a dielectric is fairly good.

Another check on our theory is to try Eq. (11.12) on atoms which have a higher frequency of excitation. For instance, it takes about 24.5 volts to pull the electron off helium, compared with the 13.5 volts required to ionize hydrogen. We would, therefore, expect that the absorption frequency ω_0 for helium would be about twice as big as for hydrogen and that α would be one-quarter as large. We expect that

$$\kappa_{\text{helium}} \approx 1.000050.$$

Experimentally,

$$\kappa_{\text{helium}} = 1.000068,$$

so you see that our rough estimates are coming out on the right track. So we have understood the dielectric constant of nonpolar gas, but only qualitatively, because we have not yet used a correct atomic theory of the motions of the atomic electrons.

11-3 Polar molecules; orientation polarization

Next we will consider a molecule which carries a permanent dipole moment p_0 —such as a water molecule. With no electric field, the individual dipoles point in random directions, so the net moment per unit volume is zero. But when an electric field is applied, two things happen: First, there is an extra dipole moment induced because of the forces on the electrons; this part gives just the same kind of electronic polarizability we found for a nonpolar molecule. For very accurate work, this effect should, of course, be included, but we will neglect it for the moment. (It can always be added in at the end.) Second, the electric field tends to line up the individual dipoles to produce a net moment per unit volume. If all the dipoles in a gas were to line up, there would be a very large polarization, but that does not happen. At ordinary temperatures and electric fields the collisions of the molecules in their thermal motion keep them from lining up very much. But there is some net alignment, and so some polarization (see Fig. 11-2). The polarization that does occur can be computed by the methods of statistical mechanics we described in Chapter 40 of Vol. I.

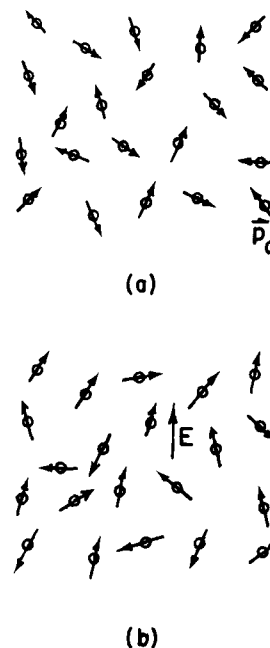


Fig. 11-2. (a) In a gas of polar molecules, the individual moments are oriented at random; the average moment in a small volume is zero. (b) When there is an electric field, there is some average alignment of the molecules.

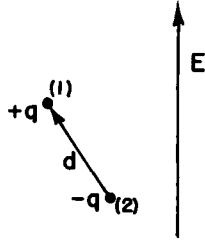


Fig. 11-3. The energy of a dipole \mathbf{p}_0 in the field \mathbf{E} is $-\mathbf{p}_0 \cdot \mathbf{E}$.

To use this method we need to know the energy of a dipole in an electric field. Consider a dipole of moment \mathbf{p}_0 in an electric field, as shown in Fig. 11-3. The energy of the positive charge is $q\phi(1)$, and the energy of the negative charge is $-q\phi(2)$. Thus the energy of the dipole is

$$U = q\phi(1) - q\phi(2) = q\mathbf{d} \cdot \nabla\phi,$$

or

$$U = -\mathbf{p}_0 \cdot \mathbf{E} = -p_0 E \cos \theta, \quad (11.14)$$

where θ is the angle between \mathbf{p}_0 and \mathbf{E} . As we would expect, the energy is lower when the dipoles are lined up with the field.

We now find out how much lining up occurs by using the methods of statistical mechanics. We found in Chapter 40 of Vol. I that in a state of thermal equilibrium, the relative number of molecules with the potential energy U is proportional to

$$e^{-U/kT}, \quad (11.15)$$

where $U(x, y, z)$ is the potential energy as a function of position. The same arguments would say that using Eq. (11.14) for the potential energy as a function of *angle*, the number of molecules at θ *per unit solid angle* is proportional to $e^{-U/kT}$.

Letting $n(\theta)$ be the number of molecules per unit solid angle at θ , we have

$$n(\theta) = n_0 e^{+p_0 E \cos \theta / kT}. \quad (11.16)$$

For normal temperatures and fields, the exponent is small, so we can approximate by expanding the exponential:

$$n(\theta) = n_0 \left(1 + \frac{p_0 E \cos \theta}{kT} \right). \quad (11.17)$$

We can find n_0 if we integrate (11.17) over all angles; the result should be just N , the total number of molecules per unit volume. The average value of $\cos \theta$ over all angles is zero, so the integral is just n_0 times the total solid angle 4π . We get

$$n_0 = \frac{N}{4\pi}. \quad (11.18)$$

We see from (11.17) that there will be more molecules oriented along the field ($\cos \theta = 1$) than against the field ($\cos \theta = -1$). So in any small volume containing many molecules there will be a net dipole moment per unit volume—that is, a polarization P . To calculate P , we want the vector sum of all the molecular moments in a unit volume. Since we know that the result is going to be in the direction of \mathbf{E} , we will just sum the components in that direction (the components at right angles to \mathbf{E} will sum to zero):

$$P = \sum_{\text{unit volume}} p_0 \cos \theta_i.$$

We can evaluate the sum by integrating over the angular distribution. The solid angle at θ is $2\pi \sin \theta d\theta$, so

$$P = \int_0^\pi n(\theta) p_0 \cos \theta 2\pi \sin \theta d\theta. \quad (11.19)$$

Substituting for $n(\theta)$ from (11.17), we have

$$P = -\frac{N}{2} \int_0^\pi \left(1 + \frac{p_0 E}{kT} \cos \theta \right) p_0 \cos \theta d(\cos \theta),$$

which is easily integrated to give

$$P = \frac{N p_0^2 E}{3kT}. \quad (11.20)$$

The polarization is proportional to the field E , so there will be normal dielectric behavior. Also, as we expect, the polarization depends inversely on the temperature, because at higher temperatures there is more disalignment by collisions. This $1/T$ dependence is called Curie's law. The permanent moment p_0 appears squared for the following reason: In a given electric field, the aligning force depends upon p_0 , and the mean moment that is produced by the lining up is again proportional to p_0 . The average induced moment is proportional to p_0^2 .

We should now try to see how well Eq. (11.20) agrees with experiment. Let's look at the case of steam. Since we don't know what p_0 is, we cannot compute P directly, but Eq. (11.20) does predict that $\kappa - 1$ should vary inversely as the temperature, and this we should check.

From (11.20) we get

$$\kappa - 1 = \frac{P}{\epsilon_0 E} = \frac{N p_0^2}{3 \epsilon_0 k T}, \quad (11.21)$$

so $\kappa - 1$ should vary in direct proportion to the density N , and inversely as the absolute temperature. The dielectric constant has been measured at several different pressures and temperatures, chosen such that the number of molecules in a unit volume remained fixed.* [Notice that if the measurements had all been taken at constant pressure, the number of molecules per unit volume would decrease linearly with increasing temperature and $\kappa - 1$ would vary as T^{-2} instead of as T^{-1} .] In Fig. 11-4 we plot the experimental observations for $\kappa - 1$ as a function of $1/T$. The dependence predicted by (11.21) is followed quite well.

There is another characteristic of the dielectric constant of polar molecules—its variation with the frequency of the applied field. Due to the moment of inertia of the molecules, it takes a certain amount of time for the heavy molecules to turn toward the direction of the field. So if we apply frequencies in the high microwave region or above, the polar contribution to the dielectric constant begins to fall away because the molecules cannot follow. In contrast to this, the electronic polarizability still remains the same up to optical frequencies, because of the smaller inertia in the electrons.

11-4 Electric fields in cavities of a dielectric

We now turn to an interesting but complicated question—the problem of the dielectric constant in dense materials. Suppose that we take liquid helium or liquid argon or some other nonpolar material. We still expect electronic polarization. But in a dense material, P can be large, so the field on an individual atom will be influenced by the polarization of the atoms in its close neighborhood. The question is, what electric field acts on the individual atom?

Imagine that the liquid is put between the plates of a condenser. If the plates are charged they will produce an electric field in the liquid. But there are also charges in the individual atoms, and the total field E is the sum of both of these effects. This true electric field varies very, very rapidly from point to point in the liquid. It is very high inside the atoms—particularly right next to the nucleus—and relatively small between the atoms. The potential difference between the plates is the line integral of this total field. If we ignore all the fine-grained variations, we can think of an *average* electric field E , which is just V/d . (This is the field we were using in the last chapter.) We should think of this field as the average over a space containing many atoms.

Now you might think that an “average” atom in an “average” location would feel this average field. But it is not that simple, as we can show by considering what happens if we imagine different-shaped holes in a dielectric. For instance, suppose that we cut a slot in a polarized dielectric, with the slot oriented parallel to the field, as shown in part (a) of Fig. 11-5. Since we know that $\nabla \times E = 0$, the line integral of E around the curve, Γ , which goes as shown in (b) of the figure, should

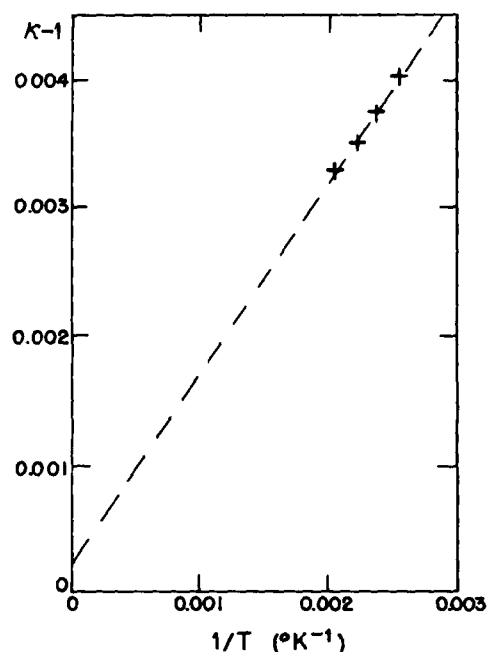


Fig. 11-4. Experimental measurements of the dielectric constant of water vapor at various temperatures.

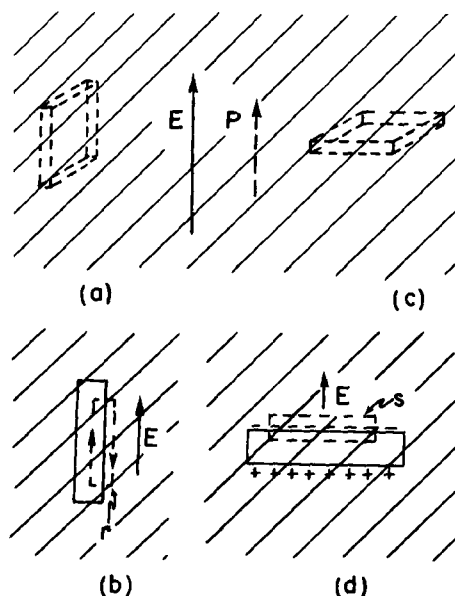


Fig. 11-5. The field in a slot cut in a dielectric depends on the shape and orientation of the slot.

* Sanger, Steiger, and Gachter, *Helvetica Physica Acta* 5, 200 (1932).

be zero. The field inside the slot must give a contribution which just cancels the part from the field outside. Therefore the field E_0 actually found in the center of a long thin slot is equal to E , the average electric field found in the dielectric.

Now consider another slot whose large sides are perpendicular to E , as shown in part (c) of Fig. 11-5. In this case, the field E_0 in the slot is not the same as E because polarization charges appear on the surfaces. If we apply Gauss' law to a surface S drawn as in (d) of the figure, we find that the field E_0 in the slot is given by

$$E_0 = E + \frac{P}{\epsilon_0}, \quad (11.22)$$

where E is again the electric field in the dielectric. (The gaussian surface contains the surface polarization charge $\sigma_{\text{pol}} = P$.) We mentioned in Chapter 10 that $\epsilon_0 E + P$ is often called D , so $\epsilon_0 E_0 = D_0$ is equal to D in the dielectric.

Earlier in the history of physics, when it was supposed to be very important to define every quantity by direct experiment, people were delighted to discover that they could define what they meant by E and D in a dielectric without having to crawl around between the atoms. The average field E is numerically equal to the field E_0 that would be measured in a slot cut parallel to the field. And the field D could be measured by finding E_0 in a slot cut normal to the field. But nobody ever measures them that way anyway, so it was just one of those philosophical things.

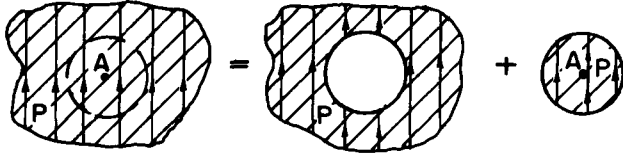


Fig. 11-6. The field at any point A in a dielectric can be considered as the sum of the field in a spherical hole plus the field due to a spherical plug.

For most liquids which are not too complicated in structure, we could expect that an atom finds itself, on the average, surrounded by the other atoms in what would be a good approximation to a *spherical hole*. And so we should ask: "What would be the field in a spherical hole?" We can find out by noticing that if we imagine carving out a spherical hole in a uniformly polarized material, we are just removing a sphere of polarized material. (We must imagine that the polarization is "frozen in" before we cut out the hole.) By superposition, however, the fields inside the dielectric, before the sphere was removed, is the sum of the fields from all charges outside the spherical volume plus the fields from the charges within the polarized sphere. That is, if we call E the field in the uniform dielectric, we can write

$$E = E_{\text{hole}} + E_{\text{plug}}, \quad (11.23)$$

where E_{hole} is the field in the hole and E_{plug} is the field inside a sphere which is uniformly polarized (see Fig. 11-6). The fields due to a uniformly polarized sphere are shown in Fig. 11-7. The electric field inside the sphere is uniform, and its value is

$$E_{\text{plug}} = -\frac{P}{3\epsilon_0}. \quad (11.24)$$

Using (11.23), we get

$$E_{\text{hole}} = E + \frac{P}{3\epsilon_0}. \quad (11.25)$$

The field in a spherical cavity is greater than the average field by the amount $P/3\epsilon_0$. (The spherical hole gives a field 1/3 of the way between a slot parallel to the field and a slot perpendicular to the field.)

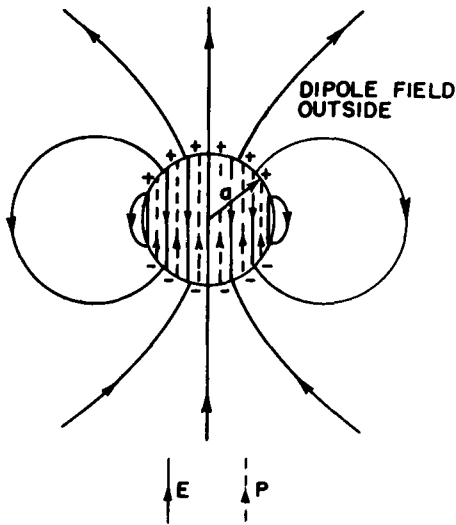


Fig. 11-7. The electric field of a uniformly polarized sphere.

11-5 The dielectric constant of liquids; the Clausius-Mossotti equation

In a liquid we expect that the field which will polarize an individual atom is more like E_{hole} than just E . If we use the E_{hole} of (11.25) for the polarizing field in

Eq. (11.6), then Eq. (11.8) becomes

$$P = N\alpha\epsilon_0\left(E + \frac{P}{3\epsilon_0}\right), \quad (11.26)$$

or

$$P = \frac{N\alpha}{1 - (N\alpha/3)} \epsilon_0 E. \quad (11.27)$$

Remembering that $\kappa - 1$ is just $P/\epsilon_0 E$, we have

$$\kappa - 1 = \frac{N\alpha}{1 - (N\alpha/3)}, \quad (11.28)$$

which gives us the dielectric constant of a liquid in terms of α , the atomic polarizability. This is called the *Clausius-Mossotti* equation.

Whenever $N\alpha$ is very small, as it is for a gas (because the density N is small), then the term $N\alpha/3$ can be neglected compared with 1, and we get our old result, Eq. (11.9), that

$$\kappa - 1 = N\alpha. \quad (11.29)$$

Let's compare Eq. (11.28) with some experimental results. It is first necessary to look at gases for which, using the measurement of κ , we can find α from Eq. (11.29). For instance, for carbon disulfide at zero degrees centigrade the dielectric constant is 1.0029, so $N\alpha$ is 0.0029. Now the density of the gas is easily worked out and the density of the liquid can be found in handbooks. At 20°C, the density of liquid CS_2 is 381 times higher than the density of the gas at 0°C. This means that N is 381 times higher in the liquid than it is in the gas so, that—if we make the approximation that the basic atomic polarizability of the carbon disulfide doesn't change when it is condensed into a liquid— $N\alpha$ in the liquid is equal to 381 times 0.0029, or 1.11. Notice that the $N\alpha/3$ term amounts to almost 0.4, so it is quite significant. With these numbers we predict a dielectric constant of 2.76, which agrees reasonably well with the observed value of 2.64.

In Table 11-1 we give some experimental data on various materials (taken from the *Handbook of Chemistry and Physics*), together with the dielectric constants calculated from Eq. (11.28) in the way just described. The agreement between observation and theory is even better for argon and oxygen than for CS_2 —and not so good for carbon tetrachloride. On the whole, the results show that Eq. (11.28) works very well.

Table 11-1
Computation of the dielectric constants of liquids
from the dielectric constant of the gas.

Substance	Gas			Liquid				
	κ (exp)	$N\alpha$	Density	Density	Ratio*	$N\alpha$	κ (predict)	κ (exp)
CS_2	1.0029	0.0029	0.00339	1.293	381	1.11	2.76	2.64
O_2	1.000523	0.000523	0.00143	1.19	832	0.435	1.509	1.507
CCl_4	1.0030	0.0030	0.00489	1.59	325	0.977	2.45	2.24
A	1.000545	0.000545	0.00178	1.44	810	0.441	1.517	1.54

* Ratio = density of liquid/density of gas.

Our derivation of Eq. (11.28) is valid only for *electronic* polarization in liquids. It is not right for a polar molecule like H_2O . If we go through the same calculations for water, we get 13.2 for $N\alpha$, which means that the dielectric constant for the liquid is *negative*, while the observed value of κ is 80. The problem has to do with the correct treatment of the permanent dipoles, and Onsager has pointed out the right way to go. We do not have the time to treat the case now, but if you are interested it is discussed in Kittel's book, *Introduction to Solid State Physics*.

11-6 Solid dielectrics

Now we turn to the solids. The first interesting fact about solids is that there can be a permanent polarization built in—which exists even without applying an electric field. An example occurs with a material like wax, which contains long molecules having a permanent dipole moment. If you melt some wax and put a strong electric field on it when it is a liquid, so that the dipole moments get partly lined up, they will stay that way when the liquid freezes. The solid material will have a permanent polarization which remains when the field is removed. Such a solid is called an *electret*.

An electret has permanent polarization charges on its surface. It is the electrical analog of a magnet. It is not as useful, though, because free charges from the air are attracted to its surfaces, eventually cancelling the polarization charges. The electret is “discharged” and there are no visible external fields.

A permanent internal polarization P is also found occurring naturally in some crystalline substances. In such crystals, each unit cell of the lattice has an identical permanent dipole moment, as drawn in Fig. 11-8. All the dipoles point in the same direction, even with no applied electric field. Many complicated crystals have, in fact, such a polarization; we do not normally notice it because the external fields are discharged, just as for the electrets.

If these internal dipole moments of a crystal are changed, however, external fields appear because there is not time for stray charges to gather and cancel the polarization charges. If the dielectric is in a condenser, free charges will be induced on the electrodes. For example, the moments can change when a dielectric is heated, because of thermal expansion. The effect is called *pyroelectricity*. Similarly, if we change the stresses in a crystal—for instance, if we bend it—again the moment may change a little bit, and a small electrical effect, called *piezoelectricity*, can be detected.

For crystals that do not have a permanent moment, one can work out a theory of the dielectric constant that involves the electronic polarizability of the atoms. It goes much the same as for liquids. Some crystals also have rotatable dipoles inside, and the rotation of these dipoles will also contribute to κ . In ionic crystals such as NaCl there is also *ionic polarizability*. The crystal consists of a checkerboard of positive and negative ions, and in an electric field the positive ions are pulled one way and the negatives the other; there is a net relative motion of the plus and minus charges, and so a volume polarization. We could estimate the magnitude of the ionic polarizability from our knowledge of the stiffness of salt crystals, but we will not go into that subject here.

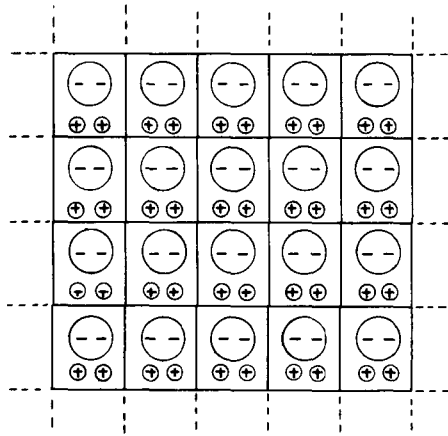


Fig. 11-8. A complex crystal lattice can have a permanent intrinsic polarization P .

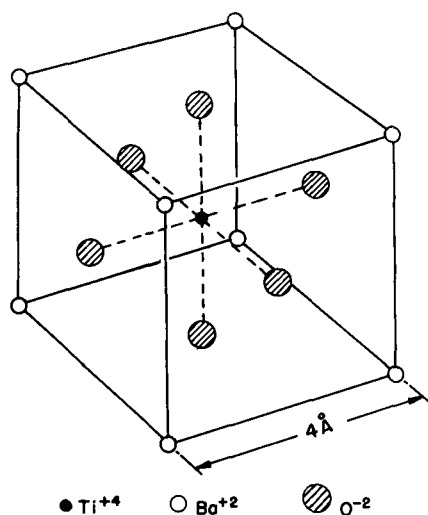


Fig. 11-9. The unit cell of BaTiO_3 . The atoms really fill up most of the space; for clarity, only the positions of their centers are shown.

11-7 Ferroelectricity; BaTiO_3

We want to describe now one special class of crystals which have, just by accident almost, a built-in permanent moment. The situation is so marginal that if we increase the temperature a little bit they lose the permanent moment completely. On the other hand, if they are nearly cubic crystals, so that their moments can be turned in different directions, we can detect a large change in the moment when an applied electric field is changed. All the moments flip over and we get a large effect. Substances which have this kind of permanent moment are called *ferroelectric*, after the corresponding ferromagnetic effects which were first discovered in iron.

We would like to explain how ferroelectricity works by describing a particular example of a ferroelectric material. There are several ways in which the ferroelectric property can originate; but we will take up only one mysterious case—that of barium titanate, BaTiO_3 . This material has a crystal lattice whose basic cell is sketched in Fig. 11-9. It turns out that above a certain temperature, specifically 118°C , barium titanate is an ordinary dielectric with an enormous dielectric constant. Below this temperature, however, it suddenly takes on a permanent moment.

In working out the polarization of solid material, we must first find what are the local fields in each unit cell. We must include the fields from the polarization

itself, just as we did for the case of a liquid. But a crystal is not a homogeneous liquid, so we cannot use for the local field what we would get in a spherical hole. If you work it out for a crystal, you find that the factor $1/3$ in Eq. (11.24) becomes slightly different, but not far from $1/3$. (For a simple cubic crystal, it is just $1/3$.) We will, therefore, assume for our preliminary discussion that the factor is $1/3$ for BaTiO_3 .

Now when we wrote Eq. (11.28) you may have wondered what would happen if $N\alpha$ became greater than 3. It appears as though κ would become negative. But that surely cannot be right. Let's see what should happen if we were gradually to increase α in a particular crystal. As α gets larger, the polarization gets bigger, making a bigger local field. But a bigger local field will polarize each atom more, raising the local fields still more. If the "give" of the atoms is enough, the process keeps going; there is a kind of feedback that causes the polarization to increase without limit—assuming that the polarization of each atom increases in proportion to the field. The "runaway" condition occurs when $N\alpha = 3$. The polarization does not become infinite, of course, because the proportionality between the induced moment and the electric field breaks down at high fields, so that our formulas are no longer correct. What happens is that the lattice gets "locked in" with a high, self-generated, internal polarization.

In the case of BaTiO_3 , there is, in addition to an electronic polarization, also a rather large ionic polarization, presumed to be due to titanium ions which can move a little within the cubic lattice. The lattice resists large motions, so after the titanium has gone a little way, it jams up and stops. But the crystal cell is then left with a permanent dipole moment.

In most crystals, this is really the situation for all temperatures that can be reached. The very interesting thing about barium titanate is that there is such a delicate condition that if $N\alpha$ is decreased just a little bit it comes unstuck. Since N decreases with increasing temperature—because of thermal expansion—we can vary $N\alpha$ by varying the temperature. Below the critical temperature it is just barely stuck, so it is easy—by applying an external field—to shift the polarization and have it lock in a different direction.

Let's see if we can analyze what happens in more detail. We call T_c the critical temperature at which $N\alpha$ is exactly 3. As the temperature increases, N goes down a little bit because of the expansion of the lattice. Since the expansion is small, we can say that near the critical temperature

$$N\alpha = 3 - \beta(T - T_c), \quad (11.30)$$

where β is a small constant, of the same order of magnitude as the thermal expansion coefficient, or about 10^{-5} to 10^{-6} per degree C. Now if we substitute this relation into Eq. (11.28), we get that

$$\kappa - 1 = \frac{3 - \beta(T - T_c)}{\beta(T - T_c)/3}.$$

Since we have assumed that $\beta(T - T_c)$ is small compared with one, we can approximate this formula by

$$\kappa - 1 = \frac{9}{\beta(T - T_c)}. \quad (11.31)$$

This relation is right, of course, only for $T > T_c$. We see that just above the critical temperature κ is enormous. Because $N\alpha$ is so close to 3, there is a tremendous magnification effect, and the dielectric constant can easily be as high as 50,000 to 100,000. It is also very sensitive to temperature. For increases in temperature, the dielectric constant goes down inversely as the temperature, but, unlike the case of a dipolar gas, for which $\kappa - 1$ goes inversely as the *absolute* temperature, for ferroelectrics it varies inversely as the difference between the absolute temperature and the critical temperature (this law is called the Curie-Weiss law).

When we lower the temperature to the critical temperature, what happens? If we imagine a lattice of unit cells like that in Fig. 11-9, we see that it is possible

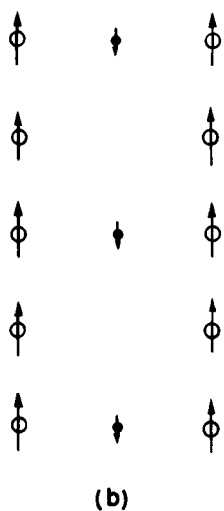
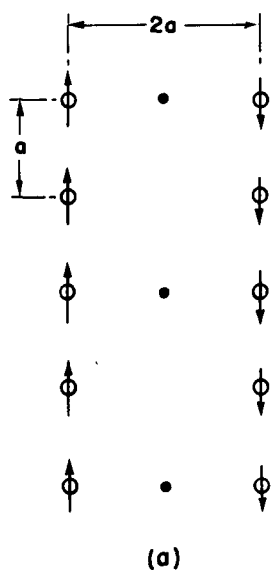


Fig. 11-10. Models of a ferroelectric: (a) corresponds to an antiferroelectric, and (b) to a normal ferroelectric.

to pick out chains of ions along vertical lines. One of them consists of alternating oxygen and titanium ions. There are other lines made up of either barium or oxygen ions, but the spacing along these lines is greater. We make a simple model to imitate this situation by imagining, as shown in Fig. 11-10(a), a series of chains of ions. Along what we call the main chain, the separation of the ions is a , which is *half* the lattice constant; the lateral distance between identical chains is $2a$. There are less-dense chains in between which we will ignore for the moment. To make the analysis a little easier, we will also suppose that all the ions on the main chain are identical. (It is not a serious simplification because all the important effects will still appear. This is one of the tricks of theoretical physics. One does a different problem because it is easier to figure out the first time—then when one understands how the thing works, it is time to put in all the complications.)

Now let's try to find out what would happen with our model. We suppose that the dipole moment of each atom is p and we wish to calculate the field at one of the atoms of the chain. We must find the sum of the fields from all the other atoms. We will first calculate the field from the dipoles in only one vertical chain; we will talk about the other chains later. The field at the distance r from a dipole in a direction along its axis is given by

$$E = \frac{1}{4\pi\epsilon_0} \frac{2p}{r^3}. \quad (11.32)$$

At any given atom, the dipoles at equal distances above and below it give fields in the same direction, so for the whole chain we get

$$E_{\text{chain}} = \frac{p}{4\pi\epsilon_0} \frac{2}{a^3} \cdot \left(2 + \frac{2}{8} + \frac{2}{27} + \frac{2}{64} + \dots \right) = \frac{p}{\epsilon_0} \frac{0.383}{a^3}. \quad (11.33)$$

It is not too hard to show that if our model were like a completely cubic crystal—that is, if the next identical lines were only the distance a away—the number 0.383 would be changed to $1/3$. In other words, if the next lines were at the distance a they would contribute only -0.050 unit to our sum. However, the next main chain we are considering is at the distance $2a$ and, as you remember from Chapter 7, the field from a periodic structure dies off exponentially with distance. Therefore these lines contribute much less than -0.050 and we can just ignore all the other chains.

It is necessary now to find out what polarizability α is needed to make the runaway process work. Suppose that the induced moment p of each atom of the chain is proportional to the field on it, as in Eq. (11.6). We get the polarizing field on the atom from E_{chain} , using Eq. (11.32). So we have the two equations

$$p = \alpha\epsilon_0 E_{\text{chain}}$$

and

$$E_{\text{chain}} = \frac{0.383}{a^3} \frac{p}{\epsilon_0}.$$

There are two solutions: E and p both zero, or

$$\alpha = \frac{a^3}{0.383},$$

with E and p both finite. Thus if α is as large as $a^3/0.383$, a permanent polarization sustained by its own field will set in. This critical equality must be reached for barium titanate at just the temperature T_c . (Notice that if α were larger than the critical value for small fields, it would decrease at larger fields and at equilibrium the same equality we have found would hold.)

For BaTiO_3 , the spacing a is 2×10^{-8} cm, so we must expect that $\alpha = 21.8 \times 10^{-24} \text{ cm}^3$. We can compare this with the known polarizabilities of the individual atoms. For oxygen, $\alpha = 30.2 \times 10^{-24} \text{ cm}^3$; we're on the right track! But for titanium, $\alpha = 2.4 \times 10^{-24} \text{ cm}^3$; rather small. To use our model we should probably take the average. (We could work out the chain again for alternating

atoms, but the result would be about the same.) So $\alpha(\text{average}) = 16.3 \times 10^{-24}$, which is not high enough to give a permanent polarization.

But wait a moment! We have so far only added up the electronic polarizabilities. There is also some ionic polarization due to the motion of the titanium ion. All we need is an ionic polarizability of $9.2 \times 10^{-24} \text{ cm}^3$. (A more precise computation using alternating atoms shows that actually 11.9×10^{-24} is needed.) To understand the properties of BaTiO_3 , we have to assume that such an ionic polarizability exists.

Why the titanium ion in barium titanate should have that much ionic polarizability is not known. Furthermore, why, at a lower temperature, it polarizes along the cube diagonal and the face diagonal equally well is not clear. If we figure out the actual size of the spheres in Fig. 11-9, and ask whether the titanium is a little bit loose in the box formed by its neighboring oxygen atoms—which is what you would hope, so that it could be easily shifted—you find quite the contrary. It fits very tightly. The *barium* atoms are slightly loose, but if you let them be the ones that move, it doesn't work out. So you see that the subject is really not one-hundred percent clear; there are still mysteries we would like to understand.

Returning to our simple model of Fig. 11-10(a), we see that the field from one chain would tend to polarize the neighboring chain in the *opposite* direction, which means that although each chain would be locked, there would be no net permanent moment per unit volume! (Although there would be no external electric effects, there are still certain thermodynamic effects one could observe.) Such systems exist, and are called antiferroelectric. So what we have explained is really an antiferroelectric. Barium titanate, however, is really like the arrangement in Fig. 11-10(b). The oxygen-titanium chains are all polarized in the same direction because there are intermediate chains of atoms in between. Although the atoms in these chains are not very polarizable, or very dense, they will be somewhat polarized, in the direction antiparallel to the oxygen-titanium chains. The small fields produced at the next oxygen-titanium chain will get it started parallel to the first. So BaTiO_3 is really ferroelectric, and it is because of the atoms in between. You may be wondering: "But what about the direct effect between the two O-Ti chains?" Remember, though, the direct effect dies off exponentially with the separation; the effect of the chain of *strong* dipoles at $2a$ can be less than the effect of a chain of weak ones at the distance a .

This completes our rather detailed report on our present understanding of the dielectric constants of gases, of liquids, and of solids.

Electrostatic Analogs

12-1 The same equations have the same solutions

The total amount of information which has been acquired about the physical world since the beginning of scientific progress is enormous, and it seems almost impossible that any one person could know a reasonable fraction of it. But it is actually quite possible for a physicist to retain a broad knowledge of the physical world rather than to become a specialist in some narrow area. The reasons for this are threefold: First, there are great principles which apply to all the different kinds of phenomena—such as the principles of the conservation of energy and of angular momentum. A thorough understanding of such principles gives an understanding of a great deal all at once. Second, there is the fact that many complicated phenomena, such as the behavior of solids under compression, really basically depend on electrical and quantum-mechanical forces, so that if one understands the fundamental laws of electricity and quantum mechanics, there is at least some possibility of understanding many of the phenomena that occur in complex situations. Finally, there is a most remarkable coincidence: *The equations for many different physical situations have exactly the same appearance.* Of course, the symbols may be different—one letter is substituted for another—but the mathematical form of the equations is the same. This means that having studied one subject, we immediately have a great deal of direct and precise knowledge about the solutions of the equations of another.

We are now finished with the subject of electrostatics, and will soon go on to study magnetism and electrodynamics. But before doing so, we would like to show that while learning electrostatics we have simultaneously learned about a large number of other subjects. We will find that the equations of electrostatics appear in several other places in physics. By a direct translation of the solutions (of course the same mathematical equations must have the same solutions) it is possible to solve problems in other fields with the same ease—or with the same difficulty—as in electrostatics.

The equations of electrostatics, we know, are

$$\nabla \cdot (\kappa E) = \frac{\rho_{\text{free}}}{\epsilon_0}, \quad (12.1)$$

$$\nabla \times E = 0. \quad (12.2)$$

(We take the equations of electrostatics with dielectrics so as to have the most general situation.) The same physics can be expressed in another mathematical form:

$$E = -\nabla\phi, \quad (12.3)$$

$$\nabla \cdot (\kappa \nabla\phi) = -\frac{\rho_{\text{free}}}{\epsilon_0}. \quad (12.4)$$

Now the point is that there are many physics problems whose mathematical equations have the same form. There is a potential (ϕ) whose gradient multiplied by a scalar function (κ) has a divergence equal to another scalar function ($-\rho/\epsilon_0$).

Whatever we know about electrostatics can immediately be carried over into that other subject, and *vice versa*. (It works both ways, of course—if the other subject has some particular characteristics that are known, then we can apply that knowledge to the corresponding electrostatic problem.) We want to consider a series of examples from different subjects that produce equations of this form.

12-1 The same equations have the same solutions

12-2 The flow of heat; a point source near an infinite plane boundary

12-3 The stretched membrane

12-4 The diffusion of neutrons; a uniform spherical source in a homogeneous medium

12-5 Irrotational fluid flow; the flow past a sphere

12-6 Illumination; the uniform lighting of a plane

12-7 The “underlying unity” of nature

12-2 The flow of heat; a point source near an infinite plane boundary

We have discussed one example earlier (Section 3-4)—the flow of heat. Imagine a block of material, which need not be homogeneous but may consist of different materials at different places, in which the temperature varies from point to point. As a consequence of these temperature variations there is a flow of heat, which can be represented by the vector \mathbf{h} . It represents the amount of heat energy which flows per unit time through a unit area perpendicular to the flow. The divergence of \mathbf{h} represents the rate per unit volume at which heat is leaving a region:

$$\nabla \cdot \mathbf{h} = \text{rate of heat out per unit volume.}$$

(We could, of course, write the equation in integral form—just as we did in electrostatics with Gauss' law—which would say that the flux through a surface is equal to the rate of change of heat energy inside the material. We will not bother to translate the equations back and forth between the differential and the integral forms, because it goes exactly the same as in electrostatics.)

The rate at which heat is generated or absorbed at various places depends, of course, on the problem. Suppose, for example, that there is a source of heat inside the material (perhaps a radioactive source, or a resistor heated by an electrical current). Let us call s the heat energy produced per unit volume per second by this source. There may also be losses (or gains) of thermal energy to other internal energies in the volume. If u is the internal energy per unit volume, $-du/dt$ will also be a "source" of heat energy. We have, then,

$$\nabla \cdot \mathbf{h} = s - \frac{du}{dt}. \quad (12.5)$$

We are not going to discuss just now the complete equation in which things change with time, because we are making an analogy to electrostatics, where nothing depends on the time. We will consider only *steady heat-flow* problems, in which constant sources have produced an equilibrium state. In these cases,

$$\nabla \cdot \mathbf{h} = s. \quad (12.6)$$

It is, of course, necessary to have another equation, which describes how the heat flows at various places. In many materials the heat current is approximately proportional to the rate of change of the temperature with position: the larger the temperature difference, the more the heat current. As we have seen, the *vector* heat current is proportional to the temperature gradient. The constant of proportionality K , a property of the material, is called the *thermal conductivity*.

$$\mathbf{h} = -K \nabla T. \quad (12.7)$$

If the properties of the material vary from place to place, then $K = K(x, y, z)$, a function of position. [Equation (12.7) is not as fundamental as (12.5), which expresses the conservation of heat energy, since the former depends upon a special property of the substance.] If now we substitute Eq. (12.7) into Eq. (12.6) we have

$$\nabla \cdot (K \nabla T) = -s, \quad (12.8)$$

which has exactly the same form as (12.4). *Steady heat-flow problems and electrostatic problems are the same.* The heat flow vector \mathbf{h} corresponds to \mathbf{E} , and the temperature T corresponds to ϕ . We have already noticed that a point heat source produces a temperature field which varies as $1/r$ and a heat flow which varies as $1/r^2$. This is nothing more than a translation of the statements from electrostatics that a point charge generates a potential which varies as $1/r$ and an electric field which varies as $1/r^2$. We can, in general, solve static heat problems as easily as we can solve electrostatic problems.

Consider a simple example. Suppose that we have a cylinder of radius a at the temperature T_1 , maintained by the generation of heat in the cylinder. (It could be, for example, a wire carrying a current, or a pipe with steam condensing inside.)

The cylinder is covered with a concentric sheath of insulating material which has a conductivity K . Say the outside radius of the insulation is b and the outside is kept at temperature T_2 (Fig. 12-1a). We want to find out at what rate heat will be lost by the wire, or steam pipe, or whatever it is in the center. Let the total amount of heat lost from a length L of the pipe be called G —which is what we are trying to find.

How can we solve this problem? We have the differential equations, but since these are the same as those of electrostatics, we have really already solved the mathematical problem. The analogous problem is that of a conductor of radius a at the potential ϕ_1 , separated from another conductor of radius b at the potential ϕ_2 , with a concentric layer of dielectric material in between, as drawn in Fig. 12-1(b). Now since the heat flow h corresponds to the electric field E , the quantity G that we want to find corresponds to the flux of the electric field from a unit length (in other words, to the electric charge per unit length over ϵ_0). We have solved the electrostatic problem by using Gauss' law. We follow the same procedure for our heat-flow problem.

From the symmetry of the situation, we know that h depends only on the distance from the center. So we enclose the pipe in a gaussian cylinder of length L and radius r . From Gauss' law, we know that the heat flow h multiplied by the area $2\pi rL$ of the surface must be equal to the total amount of heat generated inside, which is what we are calling G :

$$2\pi rLh = G \quad \text{or} \quad h = \frac{G}{2\pi rL}. \quad (12.9)$$

The heat flow is proportional to the temperature gradient:

$$h = -K \nabla T,$$

or, in this case, the magnitude of h is

$$h = -K \frac{dT}{dr}.$$

This, together with (12.9), gives

$$\frac{dT}{dr} = -\frac{G}{2\pi K L r}. \quad (12.10)$$

Integrating from $r = a$ to $r = b$, we get

$$T_2 - T_1 = -\frac{G}{2\pi K L} \ln \frac{b}{a}. \quad (12.11)$$

Solving for G , we find

$$G = \frac{2\pi K L (T_1 - T_2)}{\ln(b/a)}. \quad (12.12)$$

This result corresponds exactly to the result for the charge on a cylindrical condenser:

$$Q = \frac{2\pi\epsilon_0 L (\phi_1 - \phi_2)}{\ln(b/a)}.$$

The problems are the same, and they have the same solutions. From our knowledge of electrostatics, we also know how much heat is lost by an insulated pipe.

Let's consider another example of heat flow. Suppose we wish to know the heat flow in the neighborhood of a point source of heat located a little way beneath the surface of the earth, or near the surface of a large metal block. The localized heat source might be an atomic bomb that was set off underground, leaving an intense source of heat, or it might correspond to a small radioactive source inside a block of iron—there are numerous possibilities.

We will treat the idealized problem of a point heat source of strength G at the distance a beneath the surface of an infinite block of uniform material whose thermal conductivity is K . And we will neglect the thermal conductivity of the

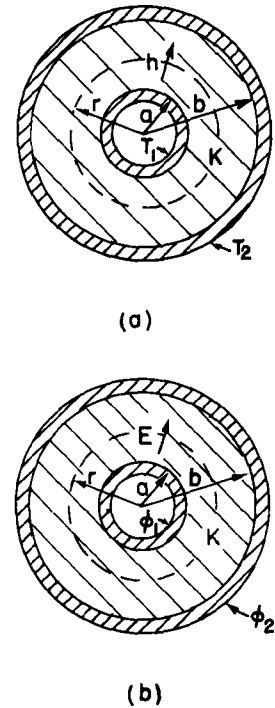


Fig. 12-1. (a) Heat flow in a cylindrical geometry. (b) The corresponding electrical problem.

air outside the material. We want to determine the distribution of the temperature on the surface of the block. How hot is it right above the source and at various places on the surface of the block?

How shall we solve it? It is like an electrostatic problem with two materials with different dielectric coefficients κ on opposite sides of a plane boundary. Aha! Perhaps it is the analog of a point charge near the boundary between a dielectric and a conductor, or something similar. Let's see what the situation is near the surface. The physical condition is that the normal component of \mathbf{h} on the surface is *zero*, since we have assumed there is no heat flow out of the block. We should ask: In what electrostatic problem do we have the condition that the normal component of the electric field \mathbf{E} (which is the analog of \mathbf{h}) is *zero* at a surface? There is none!

That is one of the things that we have to watch out for. For physical reasons, there may be certain restrictions in the kinds of mathematical conditions which arise in any one subject. So if we have analyzed the differential equation only for certain limited cases, we may have missed some kinds of solutions that can occur in other physical situations. For example, there is no material with a dielectric constant of zero, whereas a vacuum does have zero thermal conductivity. So there is no electrostatic analogy for a perfect heat insulator. We can, however, still use the same *methods*. We can try to *imagine* what would happen if the dielectric constant *were* zero. (Of course, the dielectric constant is never zero in any real situation. But we might have a case in which there is a material with a very *high* dielectric constant, so that we could neglect the dielectric constant of the air outside.)

How shall we find an electric field that has *no* component perpendicular to the surface? That is, one which is always *tangent* at the surface? You will notice that our problem is opposite to the one of a point charge near a plane conductor. There we wanted the field to be *perpendicular* to the surface, because the conductor was all at the same potential. In the electrical problem, we invented a solution by imagining a point charge behind the conducting plate. We can use the same idea again. We try to pick an "image source" that will automatically make the normal component of the field zero at the surface. The solution is shown in Fig. 12-2. An image source of the *same sign* and the same strength placed at the distance a above the surface will cause the field to be always horizontal at the surface. The normal components of the two sources cancel out.

Thus our heat flow problem is solved. The temperature everywhere is the same, by direct analogy, as the potential due to two equal point charges! The temperature T at the distance r from a single point source G in an infinite medium is

$$T = \frac{G}{4\pi K r} \quad (12.13)$$

(This, of course, is just the analog of $\phi = q/4\pi\epsilon_0 r$.) The temperature for a point source, together with its image source, is

$$T = \frac{G}{4\pi K r_1} + \frac{G}{4\pi K r_2} \quad (12.14)$$

This formula gives us the temperature everywhere in the block. Several isothermal surfaces are shown in Fig. 12-2. Also shown are lines of \mathbf{h} , which can be obtained from $\mathbf{h} = -K \nabla T$.

We originally asked for the temperature distribution on the surface. For a point on the surface at the distance ρ from the axis, $r_1 = r_2 = \sqrt{\rho^2 + a^2}$, so

$$T(\text{surface}) = \frac{1}{4\pi K} \frac{2G}{\sqrt{\rho^2 + a^2}} \quad (12.15)$$

This function is also shown in the figure. The temperature is, naturally, higher right above the source than it is farther away. This is the kind of problem that geophysicists often need to solve. We now see that it is the same kind of thing we have already been solving for electricity.

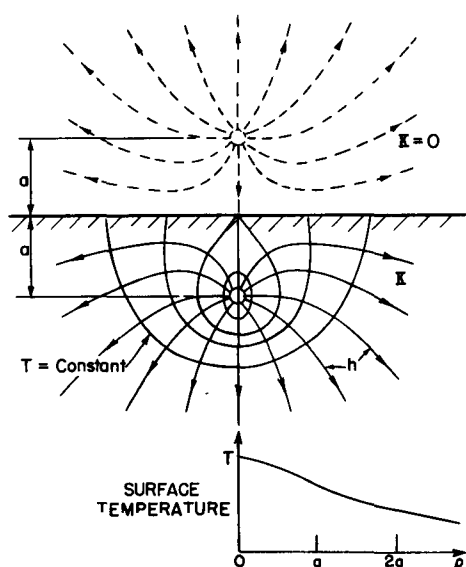


Fig. 12-2. The heat flow and isotherms near a point heat source at the distance a below the surface of a good thermal conductor. An image source is shown outside the material.

12-3 The stretched membrane

Now let us consider a completely different physical situation which, nevertheless, gives the same equations again. Consider a thin rubber sheet—a membrane—which has been stretched over a large horizontal frame (like a drumhead). Suppose now that the membrane is pushed up in one place and down in another; as shown in Fig. 12-3. Can we describe the shape of the surface? We will show how the problem can be solved when the deflections of the membrane are not too large.

There are forces in the sheet because it is stretched. If we were to make a small cut anywhere, the two sides of the cut would pull apart (see Fig. 12-4). So there is a *surface tension* in the sheet, analogous to the one-dimensional tension in a stretched string. We define the magnitude of the surface tension τ as the force *per unit length* which will just hold together the two sides of a cut such as one of those shown in Fig. 12-4.

Suppose now that we look at a vertical cross section of the membrane. It will appear as a curve, like the one in Fig. 12-5. Let u be the vertical displacement of the membrane from its normal position, and x and y the coordinates in the horizontal plane. (The cross section shown is parallel to the x -axis.)

Consider a little piece of the surface of length Δx and width Δy . There will be forces on the piece from the surface tension along each edge. The force along edge 1 of the figure will be $\tau_1 \Delta y$, directed tangent to the surface—that is, at the angle θ_1 from the horizontal. Along edge 2, the force will be $\tau_2 \Delta y$ at the angle θ_2 . (There will be similar forces on the other two edges of the piece, but we will forget them for the moment.) The net *upward* force on the piece from edges 1 and 2 is

$$\Delta F = \tau_2 \Delta y \sin \theta_2 - \tau_1 \Delta y \sin \theta_1.$$

We will limit our considerations to small distortions of the membrane, i.e., to *small slopes*: we can then replace $\sin \theta$ by $\tan \theta$, which can be written as $\partial u / \partial x$. The force is then

$$\Delta F = \left[\tau_2 \left(\frac{\partial u}{\partial x} \right)_2 - \tau_1 \left(\frac{\partial u}{\partial x} \right)_1 \right] \Delta y.$$

The quantity in brackets can be equally well written (for small Δx) as

$$\frac{\partial}{\partial x} \left(\tau \frac{\partial u}{\partial x} \right) \Delta x;$$

then

$$\Delta F = \frac{\partial}{\partial x} \left(\tau \frac{\partial u}{\partial x} \right) \Delta x \Delta y.$$

There will be another contribution to ΔF from the forces on the other two edges; the total is evidently

$$\Delta F = \left[\frac{\partial}{\partial x} \left(\tau \frac{\partial u}{\partial x} \right) + \frac{\partial}{\partial y} \left(\tau \frac{\partial u}{\partial y} \right) \right] \Delta x \Delta y. \quad (12.16)$$

The distortions of the diaphragm are caused by external forces. Let's let f represent the *upward force per unit area* on the sheet (a kind of "pressure") *from the external forces*. When the membrane is in equilibrium (the *static* case), this force must be balanced by the internal force we have just computed, Eq (12.16). That is

$$f = - \frac{\Delta F}{\Delta x \Delta y}.$$

Equation (12.16) can then be written

$$f = - \nabla \cdot (\tau \nabla u), \quad (12.17)$$

where by ∇ we now mean, of course, the two-dimensional gradient operator ($\partial/\partial x, \partial/\partial y$). We have the differential equation that relates $u(x, y)$ to the applied

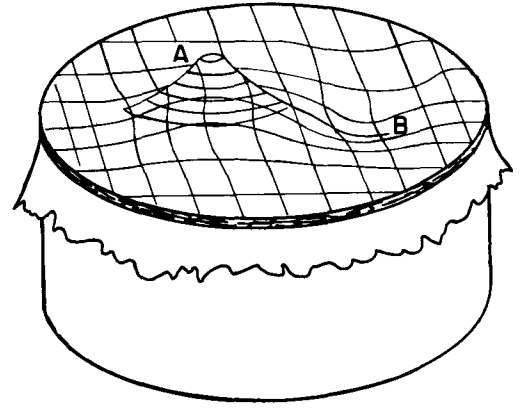


Fig. 12-3. A thin rubber sheet stretched over a cylindrical frame (like a drumhead). If the sheet is pushed up at A and down at B, what is the shape of the surface?

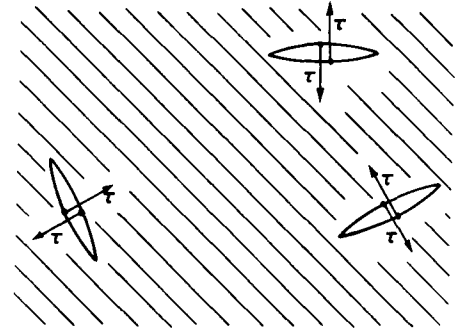


Fig. 12-4. The surface tension τ of a stretched rubber sheet is the force per unit length across a line.

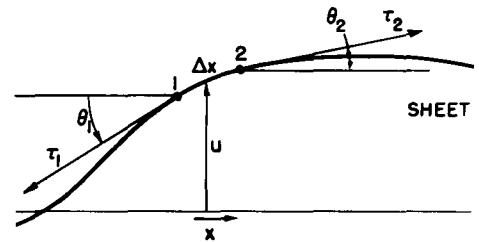


Fig. 12-5. Cross section of the deflected sheet.

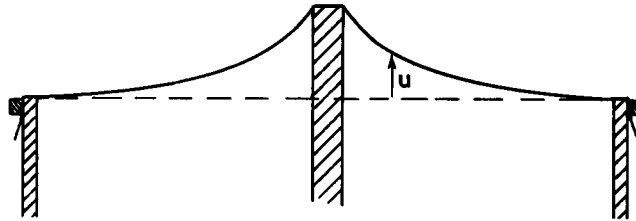
forces $f(x, y)$ and the surface tension $\tau(x, y)$, which may, in general, vary from place to place in the sheet. (The distortions of a three-dimensional elastic body are also governed by similar equations, but we will stick to two-dimensions.) We will worry only about the case in which the tension τ is constant throughout the sheet. We can then write for Eq. (12.17),

$$\nabla^2 u = -\frac{f}{\tau}. \quad (12.18)$$

We have another equation that is the same as for electrostatics!—only this time, limited to two-dimensions. The displacement u corresponds to ϕ , and f/τ corresponds to ρ/ϵ_0 . So all the work we have done for infinite plane charged sheets, or long parallel wires, or charged cylinders is directly applicable to the stretched membrane.

Suppose we push the membrane at some points up to a definite *height*—that is, we fix the value of u at some places. That is the analog of having a definite *potential* at the corresponding places in an electrical situation. So, for instance, we may make a positive “potential” by pushing up on the membrane with an object having the cross-sectional shape of the corresponding cylindrical conductor. For example, if we push the sheet up with a round rod, the surface will take on the shape shown in Fig. 12-6. The height u is the same as the electrostatic potential ϕ of a charged cylindrical rod. It falls off as $\ln(1/r)$. (The *slope*, which corresponds to the electric field E , drops off as $1/r$.)

Fig. 12-6. Cross section of a stretched rubber sheet pushed up by a round rod. The function $u(x, y)$ is the same as the electric potential $\phi(x, y)$ near a very long charged rod.



The stretched rubber sheet has often been used as a way of solving complicated *electrical* problems experimentally. The analogy is used backwards! Various rods and bars are pushed against the sheet to heights that correspond to the potentials of a set of electrodes. Measurements of the height then give the electrical potential for the electrical situation. The analogy has been carried even further. If little balls are placed on the membrane, their motion corresponds approximately to the motion of electrons in the corresponding electric field. One can actually *watch* the “electrons” move on their trajectories. This method was used to design the complicated geometry of many photomultiplier tubes (such as the ones used for scintillation counters, and the one used for controlling the headlight beams on Cadillacs). The method is still used, but the accuracy is limited. For the most accurate work, it is better to determine the fields by numerical methods, using the large electronic computing machines.

12-4 The diffusion of neutrons; a uniform spherical source in a homogeneous medium

We take another example that gives the same kind of equation, this time having to do with diffusion. In Chapter 43 of Vol. I we considered the diffusion of ions in a single gas, and of one gas through another. This time, let's take a different example—the diffusion of neutrons in a material like graphite. We choose to speak of graphite (a pure form of carbon) because carbon doesn't absorb slow neutrons. In it the neutrons are free to wander around. They travel in a straight line for several centimeters, on the average, before being scattered by a nucleus and deflected into a new direction. So if we have a large block—many meters on a side—the neutrons initially at one place will diffuse to other places. We want to find a description of their average behavior—that is, their *average flow*.

Let $N(x, y, z) \Delta V$ be the number of neutrons in the element of volume ΔV at the point (x, y, z) . Because of their motion, some neutrons will be leaving ΔV , and others will be coming in. If there are more neutrons in one region than in a nearby region, more neutrons will go from the first region to the second than come back; there will be a net flow. Following the arguments of Chapter 43 in Vol. I, we describe the flow by a flow vector \mathbf{J} . Its x -component J_x is the *net* number of neutrons that pass in unit time a unit area perpendicular to the x -direction. We found that

$$J_x = -D \frac{\partial N}{\partial x}, \quad (12.19)$$

where the diffusion constant D is given in terms of the mean velocity v , and the mean-free-path l between scatterings is given by

$$D = \frac{1}{3} lv.$$

The vector equation for \mathbf{J} is

$$\mathbf{J} = -D \nabla N. \quad (12.20)$$

The rate at which neutrons flow across any surface element da is $\mathbf{J} \cdot \mathbf{n} da$ (where, as usual, \mathbf{n} is the unit normal). The net flow *out of a volume element* is then (following the usual gaussian argument) $\nabla \cdot \mathbf{J} dV$. This flow would result in a decrease with time of the number in ΔV unless neutrons are being created in ΔV (by some nuclear process). If there are sources in the volume that generate S neutrons per unit time in a unit volume, then the net flow out of ΔV will be equal to $(S - \partial N / \partial t) \Delta V$. We have then that

$$\nabla \cdot \mathbf{J} = S - \frac{\partial N}{\partial t}. \quad (12.21)$$

Combining (12.21) with (12.20), we get the *neutron diffusion equation*

$$\nabla \cdot (-D \nabla N) = S - \frac{\partial N}{\partial t}. \quad (12.22)$$

In the static case—where $\partial N / \partial t = 0$ —we have Eq. (12.4) all over again! We can use our knowledge of electrostatics to solve problems about the diffusion of neutrons. So let's solve a problem. (You may wonder: *Why* do a problem if we have already done all the problems in electrostatics? We can do it *faster* this time because we *have* done the electrostatic problems!)

Suppose we have a block of material in which neutrons are being generated—say by uranium fission—uniformly throughout a spherical region of radius a (Fig. 12-7). We would like to know: What is the density of neutrons everywhere? How uniform is the density of neutrons in the region where they are being generated? What is the ratio of the neutron density at the center to the neutron density at the surface of the source region? Finding the answers is easy. The source density S_0 replaces the charge density ρ , so our problem is the same as the problem of a sphere of uniform charge density. Finding N is just like finding the potential ϕ . We have already worked out the fields inside and outside of a uniformly charged sphere; we can integrate them to get the potential. Outside, the potential is $Q/4\pi\epsilon_0 r$, with the total charge Q given by $4\pi a^3 \rho/3$. So

$$\phi_{\text{outside}} = \frac{\rho a^3}{3\epsilon_0 r}. \quad (12.23)$$

For points inside, the field is due only to the charge $Q(r)$ inside the sphere of radius r , $Q(r) = 4\pi r^3 \rho/3$, so

$$E = \frac{\rho r}{3\epsilon_0}. \quad (12.24)$$

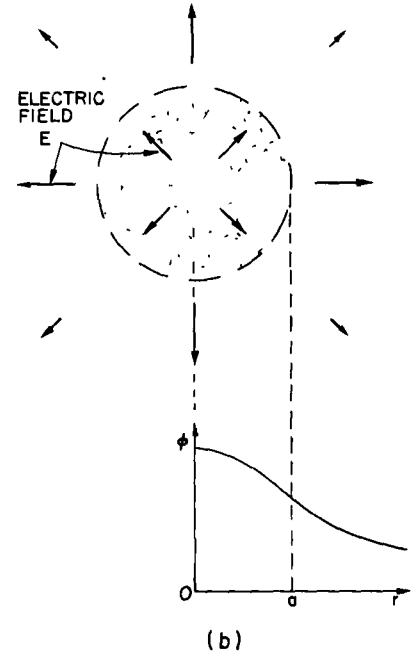
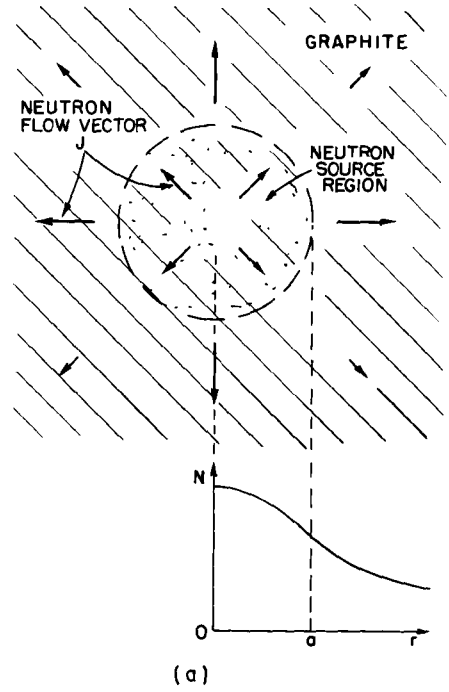


Fig. 12-7. (a) Neutrons are produced uniformly throughout a sphere of radius a in a large graphite block and diffuse outward. The neutron density N is found as a function of r , the distance from the center of the source. (b) The analogous electrostatic situation: a uniform sphere of charge, where N corresponds to ϕ and \mathbf{J} corresponds to \mathbf{E} .

The field increases linearly with r . Integrating E to get ϕ , we have

$$\phi_{\text{inside}} = -\frac{\rho r^2}{6\epsilon_0} + \text{a constant.}$$

At the radius a , ϕ_{inside} must be the same as ϕ_{outside} , so the constant must be $\rho a^2/2\epsilon_0$. (We are assuming that ϕ is zero at large distances from the source, which will correspond to N being zero for the neutrons.) Therefore,

$$\phi_{\text{inside}} = \frac{\rho}{3\epsilon_0} \left(\frac{3a^2}{2} - \frac{r^2}{2} \right). \quad (12.25)$$

We know immediately the neutron density in our other problem. The answer is

$$N_{\text{outside}} = \frac{Sa^3}{3Dr}, \quad (12.26)$$

and

$$N_{\text{inside}} = \frac{S}{3D} \left(\frac{3a^2}{2} - \frac{r^2}{2} \right). \quad (12.27)$$

N is shown as a function of r in Fig. 12-7.

Now what is the ratio of density at the center to that at the edge? At the center ($r = 0$), it is proportional to $3a^2/2$. At the edge ($r = a$) it is proportional to $2a^2/2$, so the ratio of densities is $3/2$. A uniform source doesn't produce a uniform density of neutrons. You see, our knowledge of electrostatics gives us a good start on the physics of nuclear reactors.

There are many physical circumstances in which diffusion plays a big part. The motion of ions through a liquid, or of electrons through a semiconductor, obeys the same equation. We find again and again the same equations.

12-5 Irrotational fluid flow; the flow past a sphere

Let's now consider an example which is not really a very good one, because the equations we will use will not really represent the subject with complete generality but only in an artificial idealized situation. We take up the problem of *water flow*. In the case of the stretched sheet, our equations were an approximation which was correct only for *small deflections*. For our consideration of water flow, we will not make that kind of an approximation; we must make restrictions that do not apply at all to real water. We treat only the case of the steady flow of an *incompressible, nonviscous, circulation-free* liquid. Then we represent the flow by giving the velocity $\mathbf{v}(\mathbf{r})$ as a function of position \mathbf{r} . If the motion is steady (the only case for which there is an electrostatic analog) \mathbf{v} is independent of time. If ρ is the density of the fluid, then $\rho\mathbf{v}$ is the amount of mass which passes per unit time through a unit area. By the conservation of matter, the divergence of $\rho\mathbf{v}$ will be, in general, the time rate of change of the mass of the material per unit volume. We will assume that there are no processes for the continuous creation or destruction of matter. The conservation of matter then requires that $\nabla \cdot \rho\mathbf{v} = 0$. (It should, in general, be equal to $-\partial\rho/\partial t$, but since our fluid is incompressible, ρ cannot change.) Since ρ is everywhere the same, we can factor it out, and our equation is simply

$$\nabla \cdot \mathbf{v} = 0.$$

Good! We have electrostatics again (with no charges); it's just like $\nabla \cdot \mathbf{E} = 0$. Not so! Electrostatics is *not* simply $\nabla \cdot \mathbf{E} = 0$. It is a *pair* of equations. One equation does not tell us enough; we need still an additional equation. To match electrostatics, we should have also that the *curl* of \mathbf{v} is zero. But that is not generally true for real liquids. Most liquids will ordinarily develop some circulation. So we are restricted to the situation in which there is no circulation of the fluid. Such flow is often called *irrotational*. Anyway, if we make all our assumptions, we can

imagine a case of fluid flow that is analogous to electrostatics. So we take

$$\nabla \cdot \mathbf{v} = 0 \quad (12.28)$$

and

$$\nabla \times \mathbf{v} = 0. \quad (12.29)$$

We want to emphasize that the number of circumstances in which liquid flow follows these equations is far from the great majority, but there are a few. They must be cases in which we can neglect surface tension, compressibility, and viscosity, and in which we can assume that the flow is irrotational. These assumptions are valid so rarely for real water that the mathematician John von Neumann said that people who analyze Eqs. (12.28) and (12.29) are studying “dry water”! (We take up the problem of fluid flow in more detail in Chapters 40 and 41.)

Because $\nabla \times \mathbf{v} = 0$, the velocity of “dry water” can be written as the gradient of some potential:

$$\mathbf{v} = -\nabla\psi. \quad (12.30)$$

What is the physical meaning of ψ ? There isn't any very useful meaning. The velocity can be written as the gradient of a potential simply because the flow is irrotational. And by analogy with electrostatics, ψ is called the *velocity potential*, but it is not related to a potential energy in the way that ϕ is. Since the divergence of \mathbf{v} is zero, we have

$$\nabla \cdot (\nabla\psi) = \nabla^2\psi = 0. \quad (12.31)$$

The velocity potential ψ obeys the same differential equation as the electrostatic potential in free space ($\rho = 0$).

Let's pick a problem in irrotational flow and see whether we can solve it by the methods we have learned. Consider the problem of a spherical ball falling through a liquid. If it is going too slowly, the viscous forces, which we are disregarding, will be important. If it is going too fast, little whirlpools (turbulence) will appear in its wake and there will be some circulation of the water. But if the ball is going neither too fast nor too slow, it is more or less true that the water flow will fit our assumptions, and we can describe the motion of the water by our simple equations.

It is convenient to describe what happens in a frame of reference *fixed in the sphere*. In this frame we are asking the question: How does water flow past a sphere at rest when the flow at large distances is uniform? That is, when, far from the sphere, the flow is everywhere the same. The flow near the sphere will be as shown by the streamlines drawn in Fig. 12-8. These lines, always parallel to \mathbf{v} , correspond to lines of electric field. We want to get a quantitative description for the velocity field, i.e., an expression for the velocity at any point P .

We can find the velocity from the gradient of ψ , so we first work out the potential. We want a potential that satisfies Eq. (12.31) everywhere, and which also satisfies two restrictions: (1) there is no flow in the spherical region inside the surface of the ball, and (2) the flow is constant at large distances. To satisfy (1), the component of \mathbf{v} normal to the surface of the sphere must be zero. That means that $\partial\psi/\partial r$ is zero at $r = a$. To satisfy (2), we must have $\partial\psi/\partial z = v_0$ at all points where $r \gg a$. Strictly speaking, there is no electrostatic case which corresponds exactly to our problem. It really corresponds to putting a sphere of dielectric constant *zero* in a uniform electric field. If we had worked out the solution to the problem of a sphere of a dielectric constant κ in a uniform field, then by putting $\kappa = 0$ we would immediately have the solution to this problem.

We have not actually worked out this particular electrostatic problem in detail, but let's do it now. (We could work directly on the fluid problem with \mathbf{v} and ψ , but we will use \mathbf{E} and ϕ because we are so used to them.)

The problem is: Find a solution of $\nabla^2\phi = 0$ such that $\mathbf{E} = -\nabla\phi$ is a constant, say \mathbf{E}_0 , for large r , and such that the radial component of \mathbf{E} is equal to zero at $r = a$. That is,

$$\left. \frac{\partial\phi}{\partial r} \right|_{r=a} = 0. \quad (12.32)$$

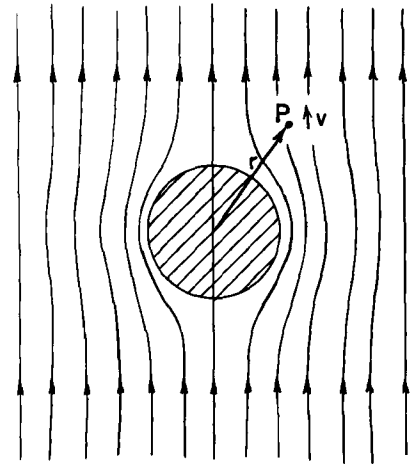


Fig. 12-8. The velocity field of irrotational fluid flow past a sphere.

Our problem involves a new kind of boundary condition, not one for which ϕ is a constant on a surface, but for which $\partial\phi/\partial r$ is a constant. That is a little different. It is not easy to get the answer immediately. First of all, without the sphere, ϕ would be $-E_0 z$. Then E would be in the z -direction and have the constant magnitude E_0 , everywhere. Now we have analyzed the case of a dielectric sphere which has a uniform polarization inside it, and we found that the field inside such a polarized sphere is a uniform field, and that outside it is the same as the field of a point dipole located at the center. So let's guess that the solution we want is a superposition of a uniform field plus the field of a dipole. The potential of a dipole (Chapter 6) is $pz/4\pi\epsilon_0 r^3$. Thus we assume that

$$\phi = -E_0 z + \frac{pz}{4\pi\epsilon_0 r^3}. \quad (12.33)$$

Since the dipole field falls off as $1/r^3$, at large distances we have just the field E_0 . Our guess will automatically satisfy condition (2) above. But what do we take for the dipole strength p ? To find out, we may use the other condition on ϕ , Eq. (12.32). We must differentiate ϕ with respect to r , but of course we must do so at a constant angle θ , so it is more convenient if we first express ϕ in terms of r and θ , rather than of z and r . Since $z = r \cos \theta$, we get

$$\phi = -E_0 r \cos \theta + \frac{p \cos \theta}{4\pi\epsilon_0 r^2}. \quad (12.34)$$

The radial component of E is

$$-\frac{\partial\phi}{\partial r} = +E_0 \cos \theta + \frac{p \cos \theta}{2\pi\epsilon_0 r^3}. \quad (12.35)$$

This must be zero at $r = a$ for all θ . This will be true if

$$p = -2\pi\epsilon_0 a^3 E_0. \quad (12.36)$$

Note carefully that if both terms in Eq. (12.35) had not had the same θ -dependence, it would not have been possible to choose p so that (12.35) turned out to be zero at $r = a$ for all angles. The fact that it works out means that we have guessed wisely in writing Eq. (12.33). Of course, when we made the guess we were looking ahead; we knew that we would need another term that (a) satisfied $\nabla^2\phi = 0$ (any real field would do that), (b) dependent on $\cos \theta$, and (c) fell to zero at large r . The dipole field is the only one that does all three.

Using (12.36), our potential is

$$\phi = -E_0 \cos \theta \left(r + \frac{a^3}{2r^2} \right). \quad (12.37)$$

The solution of the fluid flow problem can be written simply as

$$\psi = -v_0 \cos \theta \left(r + \frac{a^3}{2r^2} \right). \quad (12.38)$$

It is straightforward to find v from this potential. We will not pursue the matter further.

12-6 Illumination; the uniform lighting of a plane

In this section we turn to a completely different physical problem—we want to illustrate the great variety of possibilities. This time we will do something that leads to the same kind of *integral* that we found in electrostatics. (If we have a mathematical problem which gives us a certain integral, then we know something about the properties of that integral if it is the same integral that we had to do for another problem.) We take our example from illumination engineering. Suppose there is a light source at the distance a above a plane surface. What is the illumination of the surface? That is, what is the radiant energy per unit time arriving at a unit area of the surface? (See Fig. 12-9.) We suppose that the source is spherically

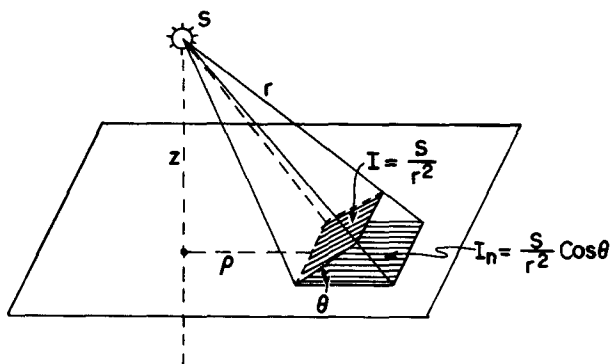


Fig. 12-9. The illumination I_n of a surface is the radiant energy per unit time arriving at a unit area of the surface.

symmetric, so that light is radiated equally in all directions. Then the amount of radiant energy which passes through a unit area *at right angles* to a light flow varies inversely as the square of the distance. It is evident that the intensity of the light in the direction normal to the flow is given by the same kind of formula as for the electric field from a point source. If the light rays meet the surface at an angle θ to the normal, then I , the energy arriving *per unit area* of the surface, is only $\cos \theta$ as great, because the same energy goes onto an area larger by $1/\cos \theta$. If we call the strength of our light source S , then I_n , the illumination of a surface, is

$$I_n = \frac{S}{r^2} \mathbf{e}_r \cdot \mathbf{n}, \quad (12.39)$$

where \mathbf{e}_r is the unit vector from the source and \mathbf{n} is the unit normal to the surface. The illumination I_n corresponds to the normal component of the electric field from a point charge of strength $4\pi\epsilon_0 S$. Knowing that, we see that for any distribution of light sources, we can find the answer by solving the corresponding electrostatic problem. We calculate the vertical component of electric field on the plane due to a distribution of charge in the same way as for that of the light sources.*

Consider the following example. We wish for some special experimental situation to arrange that the top surface of a table will have a very uniform illumination. We have available long tubular fluorescent lights which radiate uniformly along their lengths. We can illuminate the table by placing the fluorescent tubes in a regular array on the ceiling, which is at the height z above the table. What is the widest spacing b from tube to tube that we should use if we want the surface illumination to be uniform to, say, within one part in a thousand? *Answer*; (1) Find the electric field from a grid of wires with the spacing b , each charged uniformly; (2) compute the vertical component of the electric field; (3) find out what b must be so that the ripples of the field are not more than one part in a thousand.

In Chapter 7 we saw that the electric field of a grid of charged wires could be represented as a sum of terms, each one of which gave a sinusoidal variation of the field with a period of b/n , where n is an integer. The amplitude of any one of these terms is given by Eq. (7.44):

$$F_n = A_n e^{-2\pi n z/b}.$$

We need consider only $n = 1$, so long as we only want the field at points not too close to the grid. For a complete solution, we would still need to determine the coefficients A_n , which we have not yet done (although it is a straightforward calculation). Since we need only A_1 , we can estimate that its magnitude is roughly the same as that of the average field. The exponential factor would then give us directly the *relative* amplitude of the variations. If we want this factor to be 10^{-3} , we find that b must be $0.91z$. If we make the spacing of the fluorescent tubes $3/4$

* Since we are talking about *incoherent* sources whose *intensities* always add linearly, the analogous electric charges will always have the same sign. Also, our analogy applies only to the light energy arriving at the top of an opaque surface, so we must include in our integral only the sources which shine on the surface (and, naturally, not sources located below the surface!).

of the distance to the ceiling, the exponential factor is then $1/4000$, and we have a safety factor of 4, so we are fairly sure that we will have the illumination constant to one part in a thousand. (An exact calculation shows that A_1 is really twice the average field, so the exact answer is $b = 0.8z$.) It is somewhat surprising that for such a uniform illumination the allowed separation of the tubes comes out so large.

12-7 The “underlying unity” of nature

In this chapter, we wished to show that in learning electrostatics you have learned at the same time how to handle many subjects in physics, and that by keeping this in mind, it is possible to learn almost all of physics in a limited number of years.

However, a question surely suggests itself at the end of such a discussion: *Why are the equations from different phenomena so similar?* We might say: “It is the underlying unity of nature.” But what does that mean? What *could* such a statement mean? It could mean simply that the equations are similar for different phenomena; but then, of course, we have given no explanation. The “underlying unity” might mean that everything is made out of the same stuff, and therefore obeys the same equations. That sounds like a good explanation, but let us think. The electrostatic potential, the diffusion of neutrons, heat flow—are we really dealing with the same stuff? Can we really imagine that the electrostatic potential is *physically* identical to the temperature, or to the density of particles? Certainly ϕ is not *exactly the same* as the thermal energy of particles. The displacement of a membrane is certainly *not* like a temperature. Why, then, is there “an underlying unity”?

A closer look at the physics of the various subjects shows, in fact, that the equations are not really identical. The equation we found for neutron diffusion is only an approximation that is good when the distance over which we are looking is large compared with the mean free path. If we look more closely, we would see the individual neutrons running around. Certainly the motion of an individual neutron is a completely different thing from the smooth variation we get from solving the differential equation. The differential equation is an approximation, because we assume that the neutrons are smoothly distributed in *space*.

Is it possible that *this* is the clue? That the thing which is common to all the phenomena is the *space*, the framework into which the physics is put? As long as things are reasonably smooth in space, then the important things that will be involved will be the rates of change of quantities with position in space. That is why we always get an equation with a gradient. The derivatives *must* appear in the form of a gradient or a divergence; because the laws of physics are *independent of direction*, they must be expressible in vector form. The equations of electrostatics are the simplest vector equations that one can get which involve only the spatial derivatives of quantities. Any other *simple* problem—or simplification of a complicated problem—must look like electrostatics. What is common to all our problems is that they involve *space* and that we have *imitated* what is actually a complicated phenomenon by a simple differential equation.

That leads us to another interesting question. Is the same statement perhaps also true for the *electrostatic* equations? Are they also correct only as a smoothed-out imitation of a really much more complicated microscopic world? Could it be that the real world consists of little X-ons which can be seen only at *very* tiny distances? And that in our measurements we are always observing on such a large scale that we can’t see these little X-ons, and that is why we get the differential equations?

Our currently most complete theory of electrodynamics does indeed have its difficulties at very short distances. So it is possible, in principle, that these equations are smoothed-out versions of something. They appear to be correct at distances down to about 10^{-14} cm, but then they begin to look wrong. It is possible that there is some as yet undiscovered underlying “machinery,” and that the details of an underlying complexity are hidden in the smooth-looking equations—as is so

in the “smooth” diffusion of neutrons. But no one has yet formulated a successful theory that works that way.

Strangely enough, it turns out (for reasons that we do not at all understand) that the combination of relativity and quantum mechanics as we know them seems to *forbid* the invention of an equation that is fundamentally different from Eq. (12.4), and which does not at the same time lead to some kind of contradiction. Not simply a disagreement with experiment, but an *internal contradiction*. As, for example, the prediction that the sum of the probabilities of all possible occurrences is not equal to unity, or that energies may sometimes come out as complex numbers, or some other such idiocy. No one has yet made up a theory of electricity for which $\nabla^2\phi = -\rho/\epsilon_0$ is understood as a smoothed-out approximation to a mechanism underneath, and which does not lead ultimately to some kind of an absurdity. But, it must be added, it is also true that the assumption that $\nabla^2\phi = -\rho/\epsilon_0$ is valid for all distances, no matter how small, leads to absurdities of its own (the electrical energy of an electron is infinite)—absurdities from which no one yet knows an escape.

Magnetostatics

13-1 The magnetic field

The force on an electric charge depends not only on where it is, but also on how fast it is moving. Every point in space is characterized by two vector quantities which determine the force on any charge. First, there is the *electric force*, which gives a force component independent of the motion of the charge. We describe it by the electric field, E . Second, there is an additional force component, called the *magnetic force*, which depends on the velocity of the charge. This magnetic force has a strange directional character: At any particular point in space, both the *direction* of the force and its *magnitude* depend on the direction of motion of the particle: at every instant the force is always at right angles to the velocity vector; also, at any particular point, the force is always at right angles to a *fixed direction in space* (see Fig. 13-1); and finally, the magnitude of the force is proportional to the *component* of the velocity at right angles to this unique direction. It is possible to describe all of this behavior by defining the magnetic field vector B , which specifies both the unique direction in space and the constant of proportionality with the velocity, and to write the magnetic force as $qv \times B$. The total electromagnetic force on a charge can, then, be written as

$$F = q(E + v \times B). \quad (13.1)$$

This is called the *Lorentz force*.

The magnetic force is easily demonstrated by bringing a bar magnet close to a cathode-ray tube. The deflection of the electron beam shows that the presence of the magnet results in forces on the electrons transverse to their direction of motion, as we described in Chapter 12 of Vol. I.

The unit of magnetic field B is evidently one newton-second per coulomb-meter. The same unit is also one volt-second per meter². It is also called one *weber per square meter*.

13-2 Electric current; the conservation of charge

We consider first how we can understand the magnetic forces on wires carrying electric currents. In order to do this, we define what is meant by the current density. Electric currents are electrons or other charges in motion with a net drift or flow. We can represent the charge flow by a vector which gives the amount of charge passing per unit area and per unit time through a surface element at right angles to the flow (just as we did for the case of heat flow). We call this the *current density* and represent it by the vector j . It is directed along the motion of the charges. If we take a small area ΔS at a given place in the material, the amount of charge flowing across that area in a unit time is

$$j \cdot n \Delta S, \quad (13.2)$$

where n is the unit vector normal to ΔS .

The current density is related to the average flow velocity of the charges. Suppose that we have a distribution of charges whose average motion is a drift with the velocity v . As this distribution passes over a surface element ΔS , the charge Δq passing through the surface element in a time Δt is equal to the charge contained in a parallelepiped whose base is ΔS and whose height is $v \Delta t$, as shown in Fig. 13-2. The volume of the parallelepiped is the projection of ΔS at right angles to v times

13-1 The magnetic field

13-2 Electric current; the conservation of charge

13-3 The magnetic force on a current

13-4 The magnetic field of steady currents; Ampere's law

13-5 The magnetic field of a straight wire and of a solenoid; atomic currents

13-6 The relativity of magnetic and electric fields

13-7 The transformation of currents and charges

13-8 Superposition; the right-hand rule

Review: Chapter 15, Vol. I: *The Special Theory of Relativity*

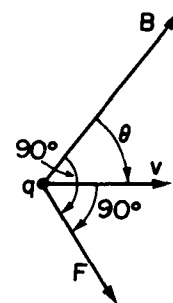


Fig. 13-1. The velocity-dependent component of the force on a moving charge is at right angles to v and to the direction of B . It is also proportional to the component of v at right angles to B , that is, to $v \sin \theta$.

$v \Delta t$, which when multiplied by the charge density ρ will give Δq . Thus

$$\Delta q = \rho v \cdot \mathbf{n} \Delta S \Delta t.$$

The charge per unit time is then $\rho v \cdot \mathbf{n} \Delta S$, from which we get

$$\mathbf{j} = \rho \mathbf{v}. \quad (13.3)$$

If the charge distribution consists of individual charges, say electrons, each with the charge q and moving with the mean velocity \mathbf{v} , then the current density is

$$\mathbf{j} = Nq\mathbf{v}, \quad (13.4)$$

where N is the number of charges per unit volume.

The total charge passing per unit time through any surface S is called the *electric current*, I . It is equal to the integral of the normal component of the flow through all of the elements of the surface:

$$I = \int_S \mathbf{j} \cdot \mathbf{n} dS \quad (13.5)$$

(see Fig. 13-3).

The current I out of a closed surface S represents the rate at which charge leaves the volume V enclosed by S . One of the basic laws of physics is that *electric charge is indestructible*; it is never lost or created. Electric charges can move from place to place but never appear from nowhere. We say that *charge is conserved*. If there is a net current out of a closed surface, the amount of charge inside must decrease by the corresponding amount (Fig. 13-4). We can, therefore, write the law of the conservation of charge as

$$\int_{\text{any closed surface}} \mathbf{j} \cdot \mathbf{n} dS = -\frac{d}{dt} (Q_{\text{inside}}). \quad (13.6)$$

The charge inside can be written as a volume integral of the charge density:

$$Q_{\text{inside}} = \int_{\text{inside } S} \rho dV. \quad (13.7)$$

If we apply (13.6) to a small volume ΔV , we know that the left-hand integral is $\nabla \cdot \mathbf{j} \Delta V$. The charge inside is $\rho \Delta V$, so the conservation of charge can also be written as

$$\nabla \cdot \mathbf{j} = -\frac{\partial \rho}{\partial t} \quad (13.8)$$

(Gauss' mathematics once again!).

13-3 The magnetic force on a current

Now we are ready to find the force on a current-carrying wire in a magnetic field. The current consists of charged particles moving with the velocity \mathbf{v} along the wire. Each charge feels a transverse force

$$\mathbf{F} = q\mathbf{v} \times \mathbf{B}$$

(Fig. 13-5a). If there are N such charges per unit volume, the number in a small volume ΔV of the wire is $N \Delta V$. The total magnetic force $\Delta \mathbf{F}$ on the volume ΔV is the sum of the forces on the individual charges, that is,

$$\Delta \mathbf{F} = (N \Delta V)(q\mathbf{v} \times \mathbf{B}).$$

But $Nq\mathbf{v}$ is just \mathbf{j} , so

$$\Delta \mathbf{F} = \mathbf{j} \times \mathbf{B} \Delta V \quad (13.9)$$

(Fig. 13-5b). The force per unit volume is $\mathbf{j} \times \mathbf{B}$.

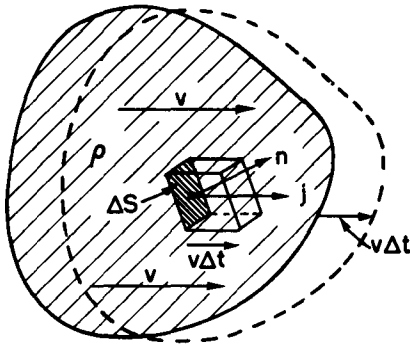


Fig. 13-2. If a charge distribution of density ρ moves with the velocity \mathbf{v} , the charge per unit time through ΔS is $\rho \mathbf{v} \cdot \mathbf{n} \Delta S$.

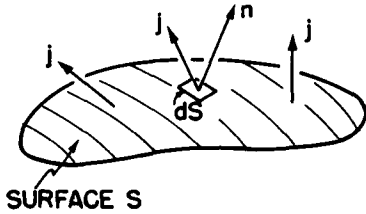


Fig. 13-3. The current I through the surface S is $\int \mathbf{j} \cdot \mathbf{n} dS$.

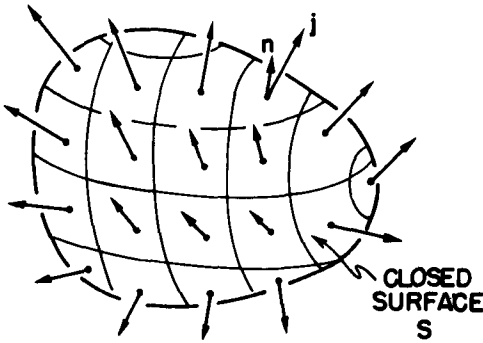


Fig. 13-4. The integral of $\mathbf{j} \cdot \mathbf{n}$ over a closed surface is the rate of change of the total charge Q inside.

If the current is uniform across a wire whose cross-sectional area is A , we may take as the volume element a cylinder with the base area A and the length ΔL . Then

$$\Delta \mathbf{F} = \mathbf{j} \times \mathbf{B} A \Delta L. \quad (13.10)$$

Now we can call $\mathbf{j}A$ the vector current \mathbf{I} in the wire. (Its magnitude is the electric current in the wire, and its direction is along the wire.) Then

$$\Delta \mathbf{F} = \mathbf{I} \times \mathbf{B} \Delta L. \quad (13.11)$$

The force per unit length on a wire is $\mathbf{I} \times \mathbf{B}$.

This equation gives the important result that the magnetic force on a wire, due to the movement of charges in it, depends only on the total current, and not on the amount of charge carried by each particle—or even its sign! The magnetic force on a wire near a magnet is easily shown by observing its deflection when a current is turned on, as was described in Chapter 1 (see Fig. 1-6).

13-4 The magnetic field of steady currents; Ampere's law

We have seen that there is a force on a wire in the presence of a magnetic field, produced, say, by a magnet. From the principle that action equals reaction we might expect that there should be a force on the source of the magnetic field, i.e., on the magnet, when there is a current through the wire.* There are indeed such forces, as is seen by the deflection of a compass needle near a current-carrying wire. Now we know that magnets feel forces from other magnets, so that means that when there is a current in a wire, the wire itself generates a magnetic field. Moving charges, then, *produce* a magnetic field. We would like now to try to discover the laws that determine how such magnetic fields are created. The question is: Given a current, what magnetic field does it make? The answer to this question was determined experimentally by three critical experiments and a brilliant theoretical argument given by Ampere. We will pass over this interesting historical development and simply say that a large number of experiments have demonstrated the validity of Maxwell's equations. We take them as our starting point. If we drop the terms involving time derivatives in these equations we get the equations of *magnetostatics*:

$$\nabla \cdot \mathbf{B} = 0 \quad (13.12)$$

and

$$c^2 \nabla \times \mathbf{B} = \frac{\mathbf{j}}{\epsilon_0}. \quad (13.13)$$

These equations are valid only if all electric charge densities are constant and all currents are steady, so that the electric and magnetic fields are not changing with time—all of the fields are “static.”

We may remark that it is rather dangerous to think that there is such a thing as a static magnetic situation, because there must be currents in order to get a magnetic field at all—and currents can come only from moving charges. “Magnetostatics” is, therefore, an approximation. It refers to a special kind of dynamic situation with *large numbers* of charges in motion, which we can approximate by a *steady* flow of charge. Only then can we speak of a current density \mathbf{j} which does not change with time. The subject should more accurately be called the study of steady currents. Assuming that all fields are steady, we drop all terms in $\partial \mathbf{E} / \partial t$ and $\partial \mathbf{B} / \partial t$ from the complete Maxwell equations, Eqs. (2.41), and obtain the two equations (13.12) and (13.13) above. Also notice that since the divergence of the curl of any vector is necessarily zero, Eq. (13.13) requires that $\nabla \cdot \mathbf{j} = 0$. This is true, by Eq. (13.8), only if $\partial \rho / \partial t$ is zero. But that must be so if \mathbf{E} is not changing with time, so our assumptions are consistent.

* We will see later, however, that such assumptions are *not* generally correct for electromagnetic forces!

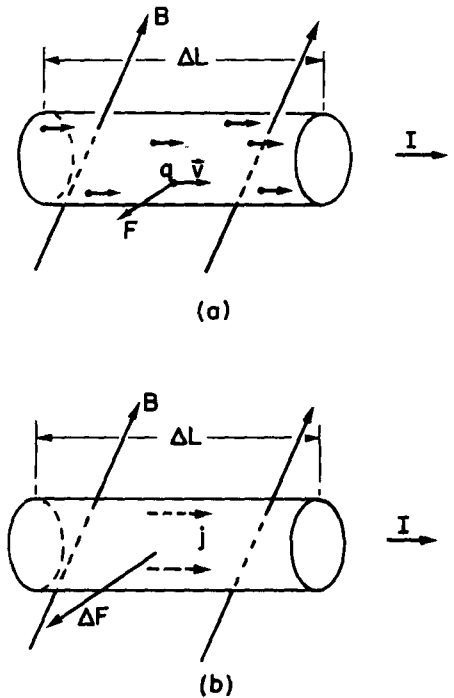


Fig. 13-5. The magnetic force on a current-carrying wire is the sum of the forces on the individual moving charges.

The requirement that $\nabla \cdot \mathbf{j} = 0$ means that we may only have charges which flow in paths that close back on themselves. They may, for instance, flow in wires that form complete loops—called circuits. The circuits may, of course, contain generators or batteries that keep the charges flowing. But they may not include condensers which are charging or discharging. (We will, of course, extend the theory later to include dynamic fields, but we want first to take the simpler case of steady currents.)

Now let us look at Eqs. (13.12) and (13.13) to see what they mean. The first one says that the divergence of \mathbf{B} is zero. Comparing it to the analogous equation in electrostatics, which says that $\nabla \cdot \mathbf{E} = \rho/\epsilon_0$, we can conclude that there is no magnetic analog of an electric charge. There are *no magnetic charges* from which lines of \mathbf{B} can emerge. If we think in terms of “lines” of the vector field \mathbf{B} , they can never start and they never stop. Then where do they come from? Magnetic fields “appear” *in the presence of currents*; they have a *curl* proportional to the current density. Wherever there are currents, there are lines of magnetic field making loops around the currents. Since lines of \mathbf{B} do not begin or end, they will often close back on themselves, making closed loops. But there can also be complicated situations in which the lines are not simple closed loops. But whatever they do, they never diverge from points. No magnetic charges have ever been discovered, so $\nabla \cdot \mathbf{B} = 0$. This much is true not only for magnetostatics, it is *always* true—even for dynamic fields.

The connection between the \mathbf{B} field and currents is contained in Eq. (13.13). Here we have a new kind of situation which is quite different from electrostatics, where we had $\nabla \times \mathbf{E} = 0$. That equation meant that the line integral of \mathbf{E} around any closed path is zero:

$$\oint_{\text{loop}} \mathbf{E} \cdot d\mathbf{s} = 0.$$

We got that result from Stokes’ theorem, which says that the integral around any closed path of *any* vector field is equal to the surface integral of the normal component of the curl of the vector (taken over any surface which has the closed loop as its periphery). Applying the same theorem to the magnetic field vector and using the symbols shown in Fig. 13-6, we get

$$\oint_{\Gamma} \mathbf{B} \cdot d\mathbf{s} = \int_S (\nabla \times \mathbf{B}) \cdot \mathbf{n} dS. \quad (13.14)$$

Taking the curl of \mathbf{B} from Eq. (13.13), we have

$$\oint_{\Gamma} \mathbf{B} \cdot d\mathbf{s} = \frac{1}{\epsilon_0 c^2} \int_S \mathbf{j} \cdot \mathbf{n} dS. \quad (13.15)$$

The integral over \mathbf{j} , according to (13.5), is the total current I through the surface S . Since for steady currents the current through S is independent of the shape of S , so long as it is bounded by the curve Γ , one usually speaks of “the current through the loop Γ .” We have, then, a general law: the circulation of \mathbf{B} around any closed curve is equal to the current I through the loop, divided by $\epsilon_0 c^2$:

$$\oint_{\Gamma} \mathbf{B} \cdot d\mathbf{s} = \frac{I_{\text{through } \Gamma}}{\epsilon_0 c^2}. \quad (13.16)$$

This law—called *Ampere’s law*—plays the same role in magnetostatics that Gauss’ law played in electrostatics. Ampere’s law alone does not determine \mathbf{B} from currents; we must, in general, also use $\nabla \cdot \mathbf{B} = 0$. But, as we will see in the next section, it can be used to find the field in special circumstances which have certain simple symmetries.

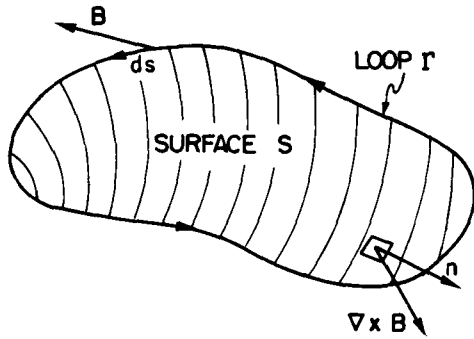


Fig. 13-6. The line integral of the tangential component of \mathbf{B} is equal to the surface integral of the normal component of $\nabla \times \mathbf{B}$.

13-5 The magnetic field of a straight wire and of a solenoid; atomic currents

We can illustrate the use of Ampere's law by finding the magnetic field near a wire. We ask: What is the field outside a long straight wire with a cylindrical cross section? We will assume something which may not be at all evident, but which is nevertheless true: that the field lines of \mathbf{B} go around the wire in closed circles. If we make this assumption, then Ampere's law, Eq. (13.16), tells us how strong the field is. From the symmetry of the problem, \mathbf{B} has the same magnitude at all points on a circle concentric with the wire (see Fig. 13-7). We can then do the line integral of $\mathbf{B} \cdot d\mathbf{s}$ quite easily; it is just the magnitude of \mathbf{B} times the circumference. If r is the radius of the circle, then

$$\oint \mathbf{B} \cdot d\mathbf{s} = B \cdot 2\pi r.$$

The total current through the loop is merely the current I in the wire, so

$$B \cdot 2\pi r = \frac{I}{\epsilon_0 c^2},$$

or

$$B = \frac{1}{4\pi\epsilon_0 c^2} \frac{2I}{r}. \quad (13.17)$$

The strength of the magnetic field drops off inversely as r , the distance from the axis of the wire. We can, if we wish, write Eq. (13.17) in vector form. Remembering that \mathbf{B} is at right angles both to \mathbf{I} and to \mathbf{r} , we have

$$\mathbf{B} = \frac{1}{4\pi\epsilon_0 c^2} \frac{2\mathbf{I} \times \mathbf{e}_r}{r}. \quad (13.18)$$

We have separated out the factor $1/4\pi\epsilon_0 c^2$, because it appears often. It is worth remembering that it is exactly 10^{-7} (in the mks system), since an equation like (13.17) is used to *define* the unit of current, the ampere. At one meter from a current of one ampere the magnetic field is 2×10^{-7} webers per square meter.

Since a current produces a magnetic field, it will exert a force on a nearby wire which is also carrying a current. In Chapter 1 we described a simple demonstration of the forces between two current-carrying wires. If the wires are parallel, each is at right angles to the \mathbf{B} field of the other; the wires should then be pushed either toward or away from each other. When currents are in the same direction, the wires attract; when the currents are moving in opposite directions, the wires repel.

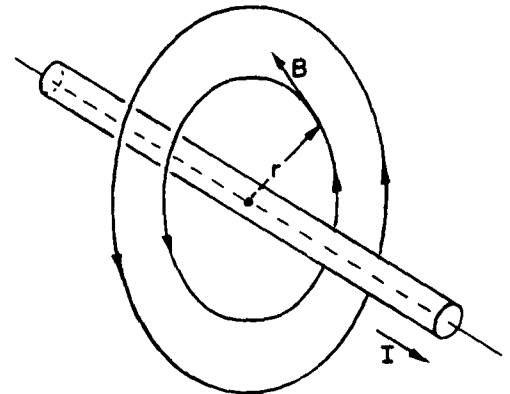


Fig. 13-7. The magnetic field outside of a long wire carrying the current I .

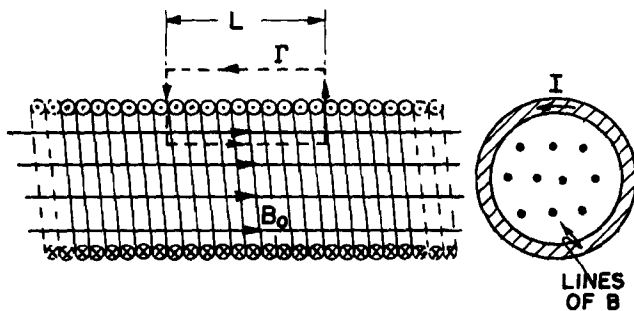


Fig. 13-8. The magnetic field of a long solenoid.

Let's take another example that can be analyzed by Ampere's law if we add some knowledge about the field. Suppose we have a long coil of wire wound in a tight spiral, as shown by the cross sections in Fig. 13-8. Such a coil is called a *solenoid*. We observe experimentally that when a solenoid is very long compared with its diameter, the field outside is very small compared with the field inside. Using just that fact, together with Ampere's law, we can find the size of the field inside.

Since the field *stays* inside (and has zero divergence), its lines must go along parallel to the axis, as shown in Fig. 13-8. That being the case, we can use Ampere's law with the rectangular "curve" Γ shown in the figure. This loop goes the distance

L inside the solenoid, where the field is, say, B_0 , then goes at right angles to the field, and returns along the outside, where the field is negligible. The line integral of \mathbf{B} for this curve is just B_0L , and it must be $1/\epsilon_0 c^2$ times the total current through Γ , which is NI if there are N turns of the solenoid in the length L . We have

$$B_0L = \frac{NI}{\epsilon_0 c^2}.$$

Or, letting n be the number of turns *per unit length* of the solenoid (that is, $n = N/L$), we get

$$B_0 = \frac{nI}{\epsilon_0 c^2}. \quad (13.19)$$

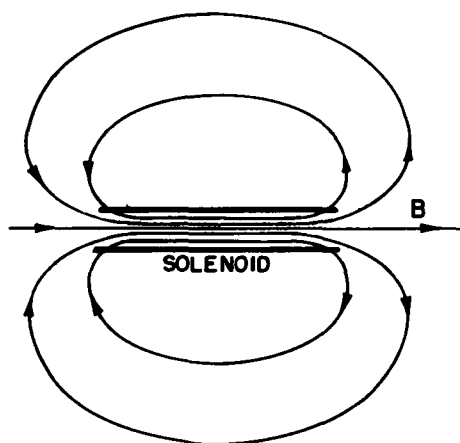


Fig. 13-9. The magnetic field outside of a solenoid.

What happens to the lines of \mathbf{B} when they get to the end of the solenoid? Presumably, they spread out in some way and return to enter the solenoid at the other end, as sketched in Fig. 13-9. Such a field is just what is observed outside of a bar magnet. But what *is* a magnet anyway? Our equations say that \mathbf{B} comes from the presence of currents. Yet we know that ordinary bars of iron (no batteries or generators) also produce magnetic fields. You might expect that there should be some other terms on the right-hand side of (13.12) or (13.13) to represent “the density of magnetic iron” or some such quantity. But there is no such term. Our theory says that the magnetic effects of iron come from some internal currents which are already taken care of by the \mathbf{j} term.

Matter is very complex when looked at from a fundamental point of view—as we saw when we tried to understand dielectrics. In order not to interrupt our present discussion, we will wait until later to deal in detail with the interior mechanisms of magnetic materials like iron. You will have to accept, for the moment, that all magnetism is produced from currents, and that in a permanent magnet there are permanent internal currents. In the case of iron, these currents come from electrons spinning around their own axes. Every electron has such a spin, which corresponds to a tiny circulating current. Of course, one electron doesn’t produce much magnetic field, but in an ordinary piece of matter there are billions and billions of electrons. Normally these spin and point every which way, so that there is no net effect. The miracle is that in a very few substances, like iron, a large fraction of the electrons spin with their axes in the same direction—for iron, two electrons from each atom takes part in this cooperative motion. In a bar magnet there are large numbers of electrons all spinning in the same direction and, as we will see, their total effect is equivalent to a current circulating on the surface of the bar. (This is quite analogous to what we found for dielectrics—that a uniformly polarized dielectric is equivalent to a distribution of charges on its surface.) It is, therefore, no accident that a bar magnet is equivalent to a solenoid.

13-6 The relativity of magnetic and electric fields

When we said that the magnetic force on a charge was proportional to its velocity, you may have wondered: “What velocity? With respect to which reference frame?” It is, in fact, clear from the definition of \mathbf{B} given at the beginning of this chapter that what this vector is will depend on what we choose as a reference frame for our specification of the velocity of charges. But we have said nothing about which is the proper frame for specifying the magnetic field.

It turns out that *any* inertial frame will do. We will also see that magnetism and electricity are not independent things—that they should always be taken together as *one* complete electromagnetic field. Although in the static case Maxwell’s equations separate into two distinct pairs, one pair for electricity and one pair for magnetism, with no apparent connection between the two fields, nevertheless, in nature itself there is a very intimate relationship between them that arises from the principle of relativity. Historically, the principle of relativity was discovered after Maxwell’s equations. It was, in fact, the study of electricity and magnetism which led ultimately to Einstein’s discovery of his principle of relativity. But let’s see

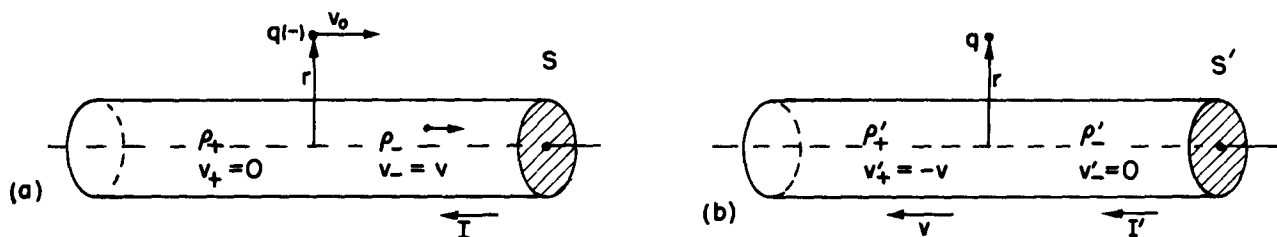


Fig. 13-10. The interaction of a current-carrying wire and a particle with the charge q as seen in two frames. In frame S (part a), the wire is at rest; in frame S' (part b), the charge is at rest.

what our knowledge of relativity would tell us about magnetic forces if we assume that the relativity principle is applicable—as it is—to electromagnetism.

Suppose we think about what happens when a negative charge moves with velocity v_0 parallel to a current-carrying wire, as in Fig. 13-10. We will try to understand what goes on in two reference frames: one fixed with respect to the wire, as in part (a) of the figure, and one fixed with respect to the particle, as in part (b). We will call the first frame S and the second S' .

In the S -frame, there is clearly a magnetic force on the particle. The force is directed toward the wire, so if the charge is moving freely we would see it curve in toward the wire. But in the S' -frame there can be no magnetic force on the particle, because its velocity is zero. Does it, therefore, stay where it is? Would we see different things happening in the two systems? The principle of relativity would say that in S' we should also see the particle move closer to the wire. We must try to understand why that would happen.

We return to our atomic description of a wire carrying a current. In a normal conductor, like copper, the electric currents come from the motion of some of the negative electrons—called the conduction electrons—while the positive nuclear charges and the remainder of the electrons stay fixed in the body of the material. We let the density of the conduction electrons be ρ_- and their velocity in S be v . The density of the charges at rest in S is ρ_+ , which must be equal to the negative of ρ_- , since we are considering an uncharged wire. There is thus no electric field outside the wire, and the force on the moving particle is just

$$\mathbf{F} = q\mathbf{v}_0 \times \mathbf{B}.$$

Using the result we found in Eq. (13.18) for the magnetic field at the distance r from the axis of a wire, we conclude that the force on the particle is directed toward the wire and has the magnitude

$$F = \frac{1}{4\pi\epsilon_0 c^2} \cdot \frac{2Iqv_0}{r}.$$

Using Eqs. (13.4) and (13.5), the current I can be written as $\rho_- vA$, where A is the area of a cross section of the wire. Then

$$F = \frac{1}{4\pi\epsilon_0 c^2} \cdot \frac{2q\rho_- Avv_0}{r}. \quad (13.20)$$

We could continue to treat the general case of arbitrary velocities for v and v_0 , but it will be just as good to look at the special case in which the velocity v_0 of the particle is the same as the velocity v of the conduction electrons. So we write $v_0 = v$, and Eq. (13.20) becomes

$$F = \frac{q}{2\pi\epsilon_0} \frac{\rho_- A}{r} \frac{v^2}{c^2}. \quad (13.21)$$

Now we turn our attention to what happens in S' , in which the particle is at rest and the wire is running past (toward the left in the figure) with the speed v . The positive charges moving with the wire will make some magnetic field B' at the particle. But the particle is now at rest, so there is no magnetic force on it! If there is any force on the particle, it must come from an electric field. It must

be that the moving wire has produced an electric field. But it can do that only if it appears *charged*—it must be that a neutral wire with a current appears to be charged when set in motion.

We must look into this. We must try to compute the charge density in the wire in S' from what we know about it in S . One might, at first, think they are the same; but we know that lengths are changed between S and S' (see Chapter 15, Vol. I), so volumes will change also. Since the charge *densities* depend on the volume occupied by charges, the densities will change, too.

Before we can decide about the charge *densities* in S' , we must know what happens to the electric *charge* of a bunch of electrons when the charges are moving. We know that the apparent mass of a particle changes by $1/\sqrt{1 - v^2/c^2}$. Does its charge do something similar? No! *Charges* are always the *same*, moving or not. Otherwise we would not always observe that the total charge is conserved.

Suppose that we take a block of material, say a conductor, which is initially uncharged. Now we heat it up. Because the electrons have a different mass than the protons, the velocities of the electrons and of the protons will change by different amounts. If the charge of a particle depended on the speed of the particle carrying it, in the heated block the charge of the electrons and protons would no longer balance. A block would become charged when heated. As we have seen earlier, a very small fractional change in the charge of all the electrons in a block would give rise to enormous electric fields. No such effect has ever been observed.

Also, we can point out that the mean speed of the electrons in matter depends on its chemical composition. If the charge on an electron changed with speed, the net charge in a piece of material would be changed in a chemical reaction. Again, a straightforward calculation shows that even a very small dependence of charge on speed would give enormous fields from the simplest chemical reactions. No such effect is observed, and we conclude that the electric charge of a single particle is independent of its state of motion.

So the charge q on a particle is an invariant scalar quantity, independent of the frame of reference. That means that in any frame the charge density of a distribution of electrons is just proportional to the number of electrons per unit volume. We need only worry about the fact that the volume *can* change because of the relativistic contraction of distances.

We now apply these ideas to our moving wire. If we take a length L_0 of the wire, in which there is a charge density ρ_0 of *stationary* charges, it will contain the total charge $Q = \rho_0 L_0 A_0$. If the same charges are observed in a different frame to be moving with velocity v , they will all be found in a piece of the material with the *shorter* length

$$L = L_0 \sqrt{1 - v^2/c^2}, \quad (13.22)$$

but with the same area A_0 (since dimensions transverse to the motion are unchanged). See Fig. 13-11.

If we call ρ the density of charges in the frame in which they are moving, the total charge Q will be $\rho L A_0$. This must also be equal to $\rho_0 L_0 A_0$, because charge is the same in any system, so that $\rho L = \rho_0 L_0$ or, from (13.22),

$$\rho = \frac{\rho_0}{\sqrt{1 - v^2/c^2}}. \quad (13.23)$$

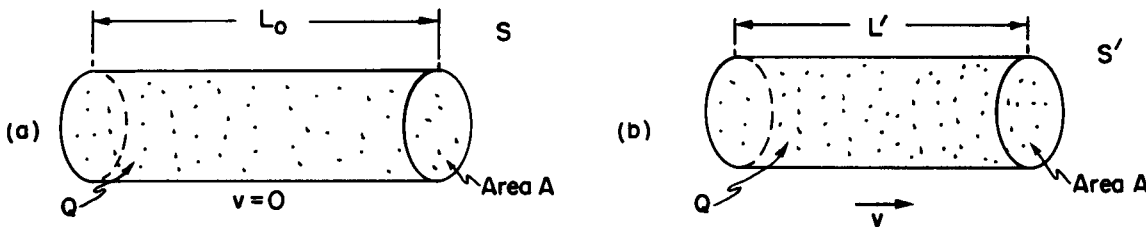


Fig. 13-11. If a distribution of charged particles at rest has the charge density ρ_0 , the same charges will have the density $\rho = \rho_0 / \sqrt{1 - v^2/c^2}$ when seen from a frame with the relative velocity v .

The charge density of a moving *distribution* of charges varies in the same way as the relativistic mass of a particle.

We now use this general result for the positive charge density ρ_+ of our wire. These charges are at rest in frame S . In S' , however, where the wire moves with the speed v , the positive charge density becomes

$$\rho'_+ = \frac{\rho_+}{\sqrt{1 - v^2/c^2}}. \quad (13.24)$$

The *negative* charges are at rest in S' . So they have their “rest density” ρ_0 in this frame. In Eq. (13.23) $\rho_0 = \rho'_-$, because they have the density ρ'_- when the wire is at rest, i.e., in frame S , where the speed of the negative charges is v . For the conduction electrons, we then have that

$$\rho_- = \frac{\rho'_-}{\sqrt{1 - v^2/c^2}}, \quad (13.25)$$

or

$$\rho'_- = \rho_- \sqrt{1 - v^2/c^2}. \quad (13.26)$$

Now we can see why there are electric fields in S' —because in this frame the wire has the net charge density ρ' given by

$$\rho' = \rho'_+ + \rho'_-.$$

Using (13.24) and (13.26), we have

$$\rho' = \frac{\rho_+}{\sqrt{1 - v^2/c^2}} + \rho_- \sqrt{1 - v^2/c^2}.$$

Since the stationary wire is neutral, $\rho_- = -\rho_+$, and we have

$$\rho' = \rho_+ \frac{v^2/c^2}{\sqrt{1 - v^2/c^2}}. \quad (13.27)$$

Our moving wire is positively charged and will produce an electric field E' at the external stationary particle. We have already solved the electrostatic problem of a uniformly charged cylinder. The electric field at the distance r from the axis of the cylinder is

$$E' = \frac{\rho' A}{2\pi\epsilon_0 r} = \frac{\rho_+ A v^2/c^2}{2\pi\epsilon_0 r \sqrt{1 - v^2/c^2}}. \quad (13.28)$$

The force on the negatively charged particle is toward the wire. We have, at least, a force in the same direction from the two points of view; the electric force in S' has the same direction as the magnetic force in S .

The magnitude of the force in S' is

$$F' = \frac{q}{2\pi\epsilon_0} \frac{\rho_+ A}{r} \frac{v^2/c^2}{\sqrt{1 - v^2/c^2}}. \quad (13.29)$$

Comparing this result for F' with our result for F in Eq. (13.21), we see that the magnitudes of the forces are almost identical from the two points of view. In fact,

$$F' = \frac{F}{\sqrt{1 - v^2/c^2}}, \quad (13.30)$$

so for the small velocities we have been considering, the two forces are equal. We can say that for low velocities, at least, we understand that magnetism and electricity are just “two ways of looking at the same thing.”

But things are even better than that. If we take into account the fact that *forces* also transform when we go from one system to the other, we find that the two ways of looking at what happens do indeed give the same *physical* result for any velocity.

One way of seeing this is to ask a question like: What transverse momentum will the particle have after the force has acted for a little while? We know from Chapter 16 of Vol. I that the transverse momentum of a particle should be the same in both the S - and S' -frames. Calling the transverse coordinate y , we want to compare Δp_y and $\Delta p'_y$. Using the relativistically correct equation of motion, $F = dp/dt$, we expect that after the time Δt our particle will have a transverse momentum Δp_y in the S -system given by

$$\Delta p_y = F \Delta t. \quad (13.31)$$

In the S' -system, the transverse momentum will be

$$\Delta p'_y = F' \Delta t'. \quad (13.32)$$

We must, of course, compare Δp_y and $\Delta p'_y$ for corresponding time intervals Δt and $\Delta t'$. We have seen in Chapter 15 of Vol. I that the time intervals referred to a *moving* particle appear to be *longer* than those in the rest system of the particle. Since our particle is initially at rest in S' , we expect, for small Δt , that

$$\Delta t = \frac{\Delta t'}{\sqrt{1 - v^2/c^2}}, \quad (13.33)$$

and everything comes out O.K. From (13.31) and (13.32),

$$\frac{\Delta p'_y}{\Delta p_y} = \frac{F' \Delta t'}{F \Delta t},$$

which is just = 1 if we combine (13.30) and (13.33).

We have found that we get the same physical result whether we analyze the motion of a particle moving along a wire in a coordinate system at rest with respect to the wire, or in a system at rest with respect to the particle. In the first instance, the force was purely “magnetic,” in the second, it was purely “electric.” The two points of view are illustrated in Fig. 13-12 (although there is still a magnetic field B' in the second frame, it produces no forces on the stationary particle).

If we had chosen still another coordinate system, we would have found a different mixture of E and B fields. Electric and magnetic forces are part of *one* physical phenomenon—the electromagnetic interactions of particles. The separation of this interaction into electric and magnetic parts depends very much on the reference frame chosen for the description. But a complete electromagnetic description is invariant; electricity and magnetism taken together are consistent with Einstein’s relativity.

Since electric and magnetic fields appear in different mixtures if we change our frame of reference, we must be careful about how we look at the fields E and B . For instance, if we think of “lines” of E or B , we must not attach too much reality to them. The lines may disappear if we try to observe them from a different coordinate system. For example, in system S' there are electric field lines, which we do *not* find “moving past us with velocity v in system S .” In system S there are no electric field lines at all! Therefore it makes no sense to say something like: When I move a magnet, it takes its field with it, so the lines of B are also moved. There is no way to make sense, in general, out of the idea of “the speed of a moving field line.” The fields are our way of describing what goes on at a point in space. In particular, E and B tell us about the forces that will act on a moving particle. The question “What is the force on a charge from a *moving* magnetic field?” doesn’t mean anything precise. The force is given by the values of E and B at the charge, and the formula (13.1) is not to be altered if the *source* of E or B is moving (it is the values of E and B that will be altered by the motion). Our mathematical description deals only with the fields as a function of x , y , z , and t *with respect to some inertial frame*.

We will later be speaking of “a *wave* of electric and magnetic fields travelling through space,” as, for instance, a light wave. But that is like speaking of a *wave* travelling on a string. We don’t then mean that some part of the *string* is moving

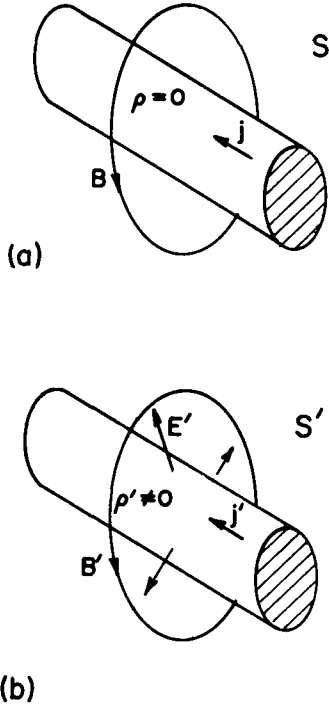


Fig. 13-12. In frame S the charge density is zero and the current density is j . There is only a magnetic field. In S' , there is a charge density ρ' , and a different current density j' . The magnetic field B' is different and there is an electric field E' .

in the direction of the wave, we mean that the *displacement* of the string appears first at one place and later at another. Similarly, in an electromagnetic wave, the *wave* travels, but the magnitude of the fields *change*. So in the future when we—or someone else—speaks of a “moving” field, you should think of it as just a handy, short way of describing a changing field in some circumstances.

13-7 The transformation of currents and charges

You may have worried about the simplification we made above when we took the same velocity v for the particle and for the conduction electrons in the wire. We could go back and carry through the analysis again for two different velocities, but it is easier to simply notice that charge and current density are the components of a four-vector (see Chapter 17, Vol. I).

We have seen that if ρ_0 is the density of the charges in their rest frame, then in a frame in which they have the velocity v , the density is

$$\rho = \frac{\rho_0}{\sqrt{1 - v^2/c^2}}.$$

In that frame their current density is

$$\mathbf{j} = \rho \mathbf{v} = \frac{\rho_0 \mathbf{v}}{\sqrt{1 - v^2/c^2}}. \quad (13.34)$$

Now we know that the energy U and momentum \mathbf{p} of a particle moving with velocity \mathbf{v} are given by

$$U = \frac{m_0 c^2}{\sqrt{1 - v^2/c^2}}, \quad \mathbf{p} = \frac{m_0 \mathbf{v}}{\sqrt{1 - v^2/c^2}},$$

where m_0 is its rest mass. We also know that U and \mathbf{p} form a relativistic four-vector. Since ρ and \mathbf{j} depend on the velocity \mathbf{v} exactly as do U and \mathbf{p} , we can conclude that ρ and \mathbf{j} are *also* the components of a relativistic four-vector. This property is the key to a general analysis of the field of a wire moving with any velocity, which we would need if we want to do the problem again with the velocity v_0 of the particle different from the velocity of the conduction electrons.

If we wish to transform ρ and \mathbf{j} to a coordinate system moving with a velocity u in the x -direction, we know that they transform just like t and (x, y, z) , so that we have (see Chapter 15, Vol. I)

$$\begin{aligned} x' &= \frac{x - ut}{\sqrt{1 - u^2/c^2}}, & j'_x &= \frac{j_x - u\rho}{\sqrt{1 - u^2/c^2}}, \\ y' &= y, & j'_y &= j_y, \\ z' &= z, & j'_z &= j_z, \\ t' &= \frac{t - ux/c^2}{\sqrt{1 - u^2/c^2}}, & \rho' &= \frac{\rho - uj_x/c^2}{\sqrt{1 - u^2/c^2}}. \end{aligned} \quad (13.35)$$

With these equations we can relate charges and currents in one frame to those in another. Taking the charges and currents in either frame, we can solve the electromagnetic problem in that frame by using our Maxwell equations. The result we obtain *for the motions of particles* will be the same no matter which frame we choose. We will return at a later time to the relativistic transformations of the electromagnetic fields.

13-8 Superposition; the right-hand rule

We will conclude this chapter by making two further points regarding the subject of magnetostatics. First, our basic equations for the magnetic field,

$$\nabla \cdot \mathbf{B} = 0, \quad \nabla \times \mathbf{B} = \mathbf{j}/c^2 \epsilon_0,$$

are linear in \mathbf{B} and \mathbf{j} . That means that the principle of superposition also applies to magnetic fields. The field produced by two different steady currents is the sum of the individual fields from each current acting alone. Our second remark concerns the right-hand rules which we have encountered (such as the right-hand rule for the magnetic field produced by a current). We have also observed that the magnetization of an iron magnet is to be understood from the spin of the electrons in the material. The direction of the magnetic field of a spinning electron is related to its spin axis by the same right-hand rule. Because \mathbf{B} is determined by a “handed” rule—involving either a cross product or a curl—it is called an *axial* vector. (Vectors whose direction in space does not depend on a reference to a right or left hand are called *polar* vectors. Displacement, velocity, force, and \mathbf{E} , for example, are polar vectors.)

Physically observable quantities in electromagnetism are *not*, however, right- (or left-) handed. Electromagnetic interactions are symmetrical under reflection (see Chapter 52, Vol. I). Whenever magnetic forces between two sets of currents are computed, the result is invariant with respect to a change in the hand convention. Our equations lead, independently of the right-hand convention, to the end result that parallel currents attract, or that currents in opposite directions repel. (Try working out the force using “left-hand rules.”) An attraction or repulsion is a polar vector. This happens because in describing any complete interaction, we use the right-hand rule twice—once to find \mathbf{B} from currents, again to find the force this \mathbf{B} produces on a second current. Using the right-hand rule twice is the same as using the left-hand rule twice. If we were to change our conventions to a left-hand system all our \mathbf{B} fields would be reversed, but all forces—or, what is perhaps more relevant, the observed accelerations of objects—would be unchanged.

Although physicists have recently found to their surprise that *all* the laws of nature are not always invariant for mirror reflections, the laws of electromagnetism do have such a basic symmetry.

The Magnetic Field in Various Situations

14-1 The vector potential

In this chapter we continue our discussion of magnetic fields associated with steady currents—the subject of magnetostatics. The magnetic field is related to electric currents by our basic equations

$$\nabla \cdot \mathbf{B} = 0, \quad (14.1)$$

$$c^2 \nabla \times \mathbf{B} = \frac{\mathbf{j}}{\epsilon_0}. \quad (14.2)$$

We want now to solve these equations mathematically in a *general* way, that is, without requiring any special symmetry or intuitive guessing. In electrostatics, we found that there was a straightforward procedure for finding the field when the positions of all electric charges are known: One simply works out the scalar potential ϕ by taking an integral over the charges—as in Eq. (4.25). Then if one wants the electric field, it is obtained from the derivatives of ϕ . We will now show that there is a corresponding procedure for finding the magnetic field \mathbf{B} if we know the current density \mathbf{j} of all moving charges.

In electrostatics we saw that (because the curl of \mathbf{E} was always zero) it was possible to represent \mathbf{E} as the gradient of a scalar field ϕ . Now the curl of \mathbf{B} is *not* always zero, so it is not possible, in general, to represent it as a gradient. However, the *divergence* of \mathbf{B} is always zero, and this means that we can always represent \mathbf{B} as the *curl* of another vector field. For, as we saw in Section 2-8, the divergence of a curl is always zero. Thus we can always relate \mathbf{B} to a field we will call \mathbf{A} by

$$\mathbf{B} = \nabla \times \mathbf{A}. \quad (14.3)$$

Or, by writing out the components,

$$\begin{aligned} B_x &= (\nabla \times \mathbf{A})_x = \frac{\partial A_z}{\partial y} - \frac{\partial A_y}{\partial z}, \\ B_y &= (\nabla \times \mathbf{A})_y = \frac{\partial A_x}{\partial z} - \frac{\partial A_z}{\partial x}, \\ B_z &= (\nabla \times \mathbf{A})_z = \frac{\partial A_y}{\partial x} - \frac{\partial A_x}{\partial y}. \end{aligned} \quad (14.4)$$

Writing $\mathbf{B} = \nabla \times \mathbf{A}$ guarantees that Eq. (14.1) is satisfied, since, necessarily,

$$\nabla \cdot \mathbf{B} = \nabla \cdot (\nabla \times \mathbf{A}) = 0.$$

The field \mathbf{A} is called the *vector potential*.

You will remember that the scalar potential ϕ was not completely specified by its definition. If we have found ϕ for some problem, we can always find another potential ϕ' that is equally good by adding a constant:

$$\phi' = \phi + C.$$

The new potential ϕ' gives the same electric fields, since the gradient ∇C is zero; ϕ' and ϕ represent the same physics.

Similarly, we can have different vector potentials \mathbf{A} which give the same magnetic fields. Again, because \mathbf{B} is obtained from \mathbf{A} by differentiation, adding a

14-1 The vector potential

14-2 The vector potential of known currents

14-3 A straight wire

14-4 A long solenoid

14-5 The field of a small loop; the magnetic dipole

14-6 The vector potential of a circuit

14-7 The law of Biot and Savart

constant to A doesn't change anything physical. But there is even more latitude for A . We can add to A any field which is the gradient of some scalar field, without changing the physics. We can show this as follows. Suppose we have an A that gives correctly the magnetic field B for some real situation, and ask in what circumstances some other new vector potential A' will give the *same* field B if substituted into (14.3). Then A and A' must have the same curl:

$$B = \nabla \times A' = \nabla \times A.$$

Therefore

$$\nabla \times A' - \nabla \times A = \nabla \times (A' - A) = 0.$$

But if the curl of a vector is zero it must be the gradient of some scalar field, say ψ , so $A' - A = \nabla\psi$. That means that if A is a satisfactory vector potential for a problem then, for any ψ at all,

$$A' = A + \nabla\psi \quad (14.5)$$

will be an equally satisfactory vector potential, leading to the same field B .

It is usually convenient to take some of the "latitude" out of A by arbitrarily placing some other condition on it (in much the same way that we found it convenient—often—to choose to make the potential ϕ zero at large distances). We can, for instance, restrict A by choosing arbitrarily what the divergence of A must be. We can always do that without affecting B . This is because although A' and A have the same curl, and give the same B , they do not need to have the same divergence. In fact, $\nabla \cdot A' = \nabla \cdot A + \nabla^2\psi$, and by a suitable choice of ψ we can make $\nabla \cdot A'$ anything we wish.

What should we choose for $\nabla \cdot A$? The choice should be made to get the greatest mathematical convenience and will depend on the problem we are doing. For *magnetostatics*, we will make the simple choice

$$\nabla \cdot A = 0. \quad (14.6)$$

(Later, when we take up electrodynamics, we will change our choice.) Our complete definition* of A is then, for the moment, $\nabla \times A = B$ and $\nabla \cdot A = 0$.

To get some experience with the vector potential, let's look first at what it is for a uniform magnetic field B_0 . Taking our z-axis in the direction of B_0 , we must have

$$\begin{aligned} B_x &= \frac{\partial A_z}{\partial y} - \frac{\partial A_y}{\partial z} = 0, \\ B_y &= \frac{\partial A_x}{\partial z} - \frac{\partial A_z}{\partial x} = 0, \\ B_z &= \frac{\partial A_y}{\partial x} - \frac{\partial A_x}{\partial y} = B_0. \end{aligned} \quad (14.7)$$

By inspection, we see that one *possible* solution of these equations is

$$A_y = xB_0, \quad A_x = 0, \quad A_z = 0.$$

Or we could equally well take

$$A_x = -yB_0, \quad A_y = 0, \quad A_z = 0.$$

Still another solution is a linear combination of the two:

$$A_x = -\frac{1}{2}yB_0, \quad A_y = \frac{1}{2}xB_0, \quad A_z = 0. \quad (14.8)$$

* Our definition still does not uniquely determine A . For a *unique* specification we would also have to say something about how the field A behaves on some boundary, or at large distances. It is sometimes convenient, for example, to choose a field which goes to zero at large distances.

It is clear that for any particular field \mathbf{B} , the vector potential \mathbf{A} is not unique; there are many possibilities.

The third solution, Eq. (14.8), has some interesting properties. Since the x -component is proportional to $-y$ and the y -component is proportional to $+x$, \mathbf{A} must be at right angles to the vector from the z -axis, which we will call \mathbf{r}' (the "prime" is to remind us that it is *not* the vector displacement from the origin). Also, the magnitude of \mathbf{A} is proportional to $\sqrt{x^2 + y^2}$ and, hence, to r' . So \mathbf{A} can be simply written (for our uniform field) as

$$\mathbf{A} = \frac{1}{2}\mathbf{B} \times \mathbf{r}'. \quad (14.9)$$

The vector potential \mathbf{A} has the magnitude $Br'/2$ and rotates about the z -axis as shown in Fig. 14-1. If, for example, the \mathbf{B} field is the axial field inside a solenoid, then the vector potential circulates in the same sense as do the currents of the solenoid.

The vector potential for a uniform field can be obtained in another way. The circulation of \mathbf{A} on any closed loop Γ can be related to the surface integral of $\nabla \times \mathbf{A}$ by Stokes' theorem, Eq. (3.38):

$$\oint_{\Gamma} \mathbf{A} \cdot d\mathbf{s} = \int_{\text{inside } \Gamma} (\nabla \times \mathbf{A}) \cdot \mathbf{n} \, da. \quad (14.10)$$

But the integral on the right is equal to the flux of \mathbf{B} through the loop, so

$$\oint_{\Gamma} \mathbf{A} \cdot d\mathbf{s} = \int_{\text{inside } \Gamma} \mathbf{B} \cdot \mathbf{n} \, da. \quad (14.11)$$

So the circulation of \mathbf{A} around *any* loop is equal to the flux of \mathbf{B} through the loop. If we take a circular loop, of radius r' in a plane perpendicular to a uniform field \mathbf{B} , the flux is just

$$\pi r'^2 B.$$

If we choose our origin on an axis of symmetry, so that we can take \mathbf{A} as circumferential and a function only of r' , the circulation will be

$$\oint \mathbf{A} \cdot d\mathbf{s} = 2\pi r' A = \pi r'^2 B.$$

We get, as before,

$$A = \frac{Br'}{2}.$$

In the example we have just given, we have calculated the vector potential from the magnetic field, which is opposite to what one normally does. In complicated problems it is usually easier to solve for the vector potential, and then determine the magnetic field from it. We will now show how this can be done.

14-2 The vector potential of known currents

Since \mathbf{B} is determined by currents, so also is \mathbf{A} . We want now to find \mathbf{A} in terms of the currents. We start with our basic equation (14.2):

$$c^2 \nabla \times \mathbf{B} = \frac{\mathbf{j}}{\epsilon_0},$$

which means, of course, that

$$c^2 \nabla \times (\nabla \times \mathbf{A}) = \frac{\mathbf{j}}{\epsilon_0}. \quad (14.12)$$

This equation is for magnetostatics what the equation

$$\nabla \cdot \nabla \phi = -\frac{\rho}{\epsilon_0} \quad (14.13)$$

was for electrostatics.

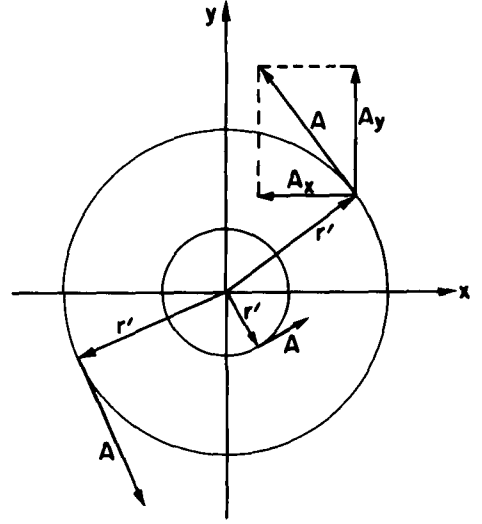


Fig. 14-1. A uniform magnetic field \mathbf{B} in the z -direction corresponds to a vector potential \mathbf{A} that rotates about the z -axis, with the magnitude $A = Br'/2$ (r' is the displacement from the z -axis).

Our equation (14.12) for the vector potential looks even more like that for ϕ if we rewrite $\nabla \times (\nabla \times \mathbf{A})$ using the vector identity Eq. (2.58):

$$\nabla \times (\nabla \times \mathbf{A}) = \nabla(\nabla \cdot \mathbf{A}) - \nabla^2 \mathbf{A}. \quad (14.14)$$

Since we have chosen to make $\nabla \cdot \mathbf{A} = 0$ (and now you see why), Eq. (14.12) becomes

$$\nabla^2 \mathbf{A} = -\frac{\mathbf{j}}{\epsilon_0 c^2}. \quad (14.15)$$

This vector equation means, of course, three equations:

$$\nabla^2 A_x = -\frac{j_x}{\epsilon_0 c^2}, \quad \nabla^2 A_y = -\frac{j_y}{\epsilon_0 c^2}, \quad \nabla^2 A_z = -\frac{j_z}{\epsilon_0 c^2}. \quad (14.16)$$

And each of these equations is *mathematically identical* to

$$\nabla^2 \phi = -\frac{\rho}{\epsilon_0}. \quad (14.17)$$

All we have learned about solving for potentials when ρ is known can be used for solving for each component of \mathbf{A} when \mathbf{j} is known!

We have seen in Chapter 4 that a general solution for the electrostatic equation (14.17) is

$$\phi(1) = \frac{1}{4\pi\epsilon_0} \int \frac{\rho(2) dV_2}{r_{12}}.$$

So we know immediately that a general solution for A_x is

$$A_x(1) = \frac{1}{4\pi\epsilon_0 c^2} \int \frac{j_x(2) dV_2}{r_{12}}, \quad (14.18)$$

and similarly for A_y and A_z . (Figure 14-2 will remind you of our conventions for r_{12} and dV_2 .) We can combine the three solutions in the vector form

$$\mathbf{A}(1) = \frac{1}{4\pi\epsilon_0 c^2} \int \frac{\mathbf{j}(2) dV_2}{r_{12}}. \quad (14.19)$$

(You can verify if you wish, by direct differentiation of components, that this integral for \mathbf{A} satisfies $\nabla \cdot \mathbf{A} = 0$ so long as $\nabla \cdot \mathbf{j} = 0$, which, as we saw, must happen for steady currents.)

We have, then, a general method for finding the magnetic field of steady currents. The principle is: the x -component of vector potential arising from a current density \mathbf{j} is the same as the electric potential ϕ that would be produced by a charge density ρ equal to j_x/c^2 —and similarly for the y - and z -components. (This principle works only with components in fixed directions. The “radial” component of \mathbf{A} does not come in the same way from the “radial” component of \mathbf{j} , for example.) So from the vector current density \mathbf{j} , we can find \mathbf{A} using Eq. (14.19)—that is, we find each component of \mathbf{A} by solving three imaginary electrostatic problems for the charge distributions $\rho_1 = j_x/c^2$, $\rho_2 = j_y/c^2$, and $\rho_3 = j_z/c^2$. Then we get \mathbf{B} by taking various derivatives of \mathbf{A} to obtain $\nabla \times \mathbf{A}$. It’s a little more complicated than electrostatics, but the same idea. We will now illustrate the theory by solving for the vector potential in a few special cases.

14-3 A straight wire

For our first example, we will again find the field of a straight wire—which we solved in the last chapter by using Eq. (14.2) and some arguments of symmetry. We take a long straight wire of radius a , carrying the steady current I . Unlike the charge on a conductor in the electrostatic case, a steady current in a wire is uniformly distributed throughout the cross section of the wire. If we choose our

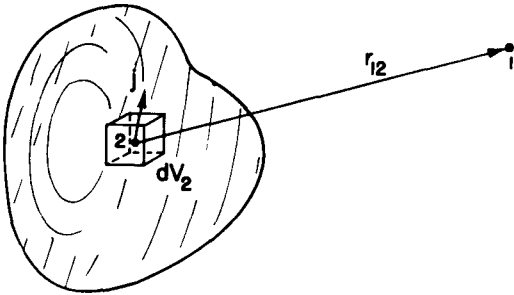


Fig. 14-2. The vector potential \mathbf{A} at point 1 is given by an integral over the current elements $\mathbf{j} dV$ at all points 2.

coordinates as shown in Fig. 14-3, the current density vector j has only a z -component. Its magnitude is

$$j_z = \frac{I}{\pi a^2} \quad (14.20)$$

inside the wire, and zero outside.

Since j_x and j_y are both zero, we have immediately

$$A_x = 0, \quad A_y = 0.$$

To get A_z we can use our solution for the electrostatic potential ϕ of a wire with a uniform charge density $\rho = j_z/c^2$. For points outside an infinite charged cylinder, the electrostatic potential is

$$\phi = -\frac{\lambda}{2\pi\epsilon_0} \ln r',$$

where $r' = \sqrt{x^2 + y^2}$ and λ is the charge per unit length, $\pi a^2 \rho$. So A_z must be

$$A_z = -\frac{\pi a^2 j_z}{2\pi\epsilon_0 c^2} \ln r'$$

for points outside a long wire carrying a uniform current. Since $\pi a^2 j_z = I$, we can also write

$$A_z = -\frac{I}{2\pi\epsilon_0 c^2} \ln r'. \quad (14.21)$$

Now we can find B from (14.4). There are only two of the six derivatives that are not zero. We get

$$B_x = -\frac{I}{2\pi\epsilon_0 c^2} \frac{\partial}{\partial y} \ln r' = -\frac{I}{2\pi\epsilon_0 c^2} \frac{y}{r'^2}, \quad (14.22)$$

$$B_y = \frac{I}{2\pi\epsilon_0 c^2} \frac{\partial}{\partial x} \ln r' = \frac{I}{2\pi\epsilon_0 c^2} \frac{x}{r'^2}, \quad (14.23)$$

$$B_z = 0.$$

We get the same result as before: B circles around the wire, and has the magnitude

$$B = \frac{1}{4\pi\epsilon_0 c^2} \frac{2I}{r'}. \quad (14.24)$$

14-4 A long solenoid

Next, we consider again the infinitely long solenoid with a circumferential current on the surface of nI per unit length. (We imagine there are n turns of wire per unit length, carrying the current I , and we neglect the slight pitch of the winding.)

Just as we have defined a "surface charge density" σ , we define here a "surface current density" J equal to the current per unit length on the surface of the solenoid (which is, of course, just the average j times the thickness of the thin winding). The magnitude of J is, here, nI . This surface current (see Fig. 14-4) has the components.

$$J_x = -J \sin \phi, \quad J_y = J \cos \phi, \quad J_z = 0.$$

Now we must find A for such a current distribution.

First, we wish to find A_z for points outside the solenoid. The result is the same as the electrostatic potential outside a cylinder with a surface charge

$$\sigma = \sigma_0 \sin \phi,$$

with $\sigma_0 = J/c^2$. We have not solved such a charge distribution, but we have done something similar. This charge distribution is equivalent to two *solid* cylinders of charge, one positive and one negative, with a slight relative displacement of their

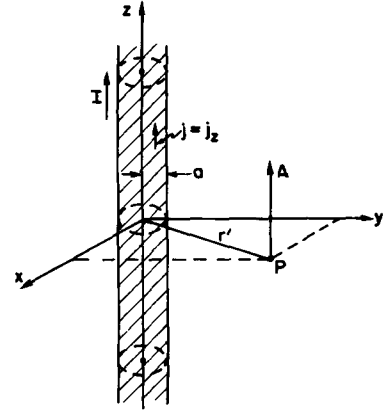


Fig. 14-3. A long cylindrical wire along the z -axis with a uniform current density j .

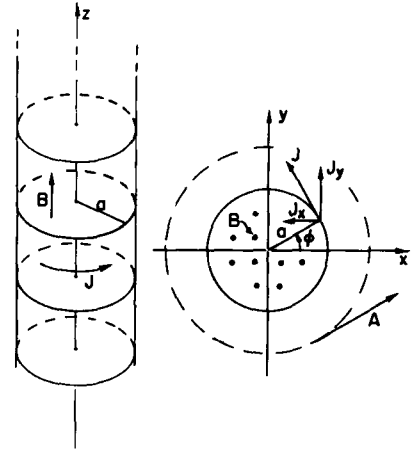


Fig. 14-4. A long solenoid with a surface current density J .

axes in the y -direction. The potential of such a pair of cylinders is proportional to the derivative with respect to y of the potential of a single uniformly charged cylinder. We could work out the constant of proportionality, but let's not worry about it for the moment.

The potential of a cylinder of charge is proportional to $\ln r'$; the potential of the pair is then

$$\phi \propto \frac{\partial \ln r'}{\partial y} = \frac{y}{r'^2}.$$

So we know that

$$A_x = -K \frac{y}{r'^2}, \quad (14.25)$$

where K is some constant. Following the same argument, we would find

$$A_y = K \frac{x}{r'^2}. \quad (14.26)$$

Although we said before that there was no *magnetic* field outside a solenoid, we find now that there *is* an *A*-field which circulates around the z -axis, as in Fig. 14-4. The question is: Is its curl zero?

Clearly, B_x and B_y are zero, and

$$\begin{aligned} B_z &= \frac{\partial}{\partial x} \left(K \frac{x}{r'^2} \right) - \frac{\partial}{\partial y} \left(-K \frac{y}{r'^2} \right) \\ &= K \left(\frac{1}{r'^2} - \frac{2x^2}{r'^4} + \frac{1}{r'^2} - \frac{2y^2}{r'^4} \right) = 0. \end{aligned}$$

So the magnetic field outside a very long solenoid is indeed zero, even though the vector potential is not.

We can check our result against something else we know: The circulation of the vector potential around the solenoid should be equal to the flux of B inside the coil (Eq. 14.11). The circulation is $A \cdot 2\pi r'$ or, since $A = K/r'$, the circulation is $2\pi K$. Notice that it is independent of r' . That is just as it should be if there is no B outside, because the flux is just the magnitude of B inside the solenoid times πa^2 . It is the same for all circles of radius $r' > a$. We have found in the last chapter that the field inside is $nI/\epsilon_0 c^2$, so we can determine the constant K :

$$2\pi K = \pi a^2 \frac{nI}{\epsilon_0 c^2},$$

or

$$K = \frac{nIa^2}{2\epsilon_0 c^2}.$$

So the vector potential *outside* has the magnitude

$$A = \frac{nIa^2}{2\epsilon_0 c^2} \frac{1}{r'}, \quad (14.27)$$

and is always perpendicular to the vector r' .

We have been thinking of a solenoidal coil of wire, but we would produce the same fields if we rotated a long cylinder with an electrostatic charge on the surface. If we have a thin cylindrical shell of radius a with a surface charge σ , rotating the cylinder makes a surface current $J = \sigma v$, where $v = a\omega$ is the velocity of the surface charge. There will then be a magnetic field $B = \sigma a \omega / \epsilon_0 c^2$ inside the cylinder.

Now we can raise an interesting question. Suppose we put a short piece of wire W perpendicular to the axis of the cylinder, extending from the axis out to the surface, and fastened to the cylinder so that it rotates with it, as in Fig. 14-5. This wire is moving in a magnetic field, so the $v \times B$ forces will cause the ends of the wire to be charged (they will charge up until the E -field from the charges just balances the $v \times B$ force). If the cylinder has a positive charge, the end of the wire at the axis will have a negative charge. By measuring the charge on the end of the

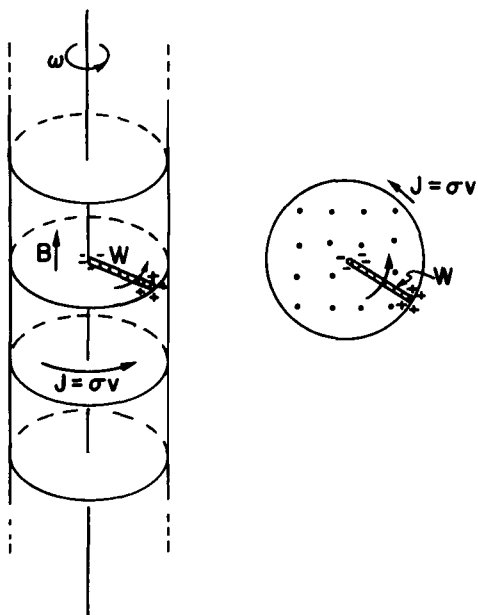


Fig. 14-5. A rotating charged cylinder produces a magnetic field inside. A short radial wire rotating with the cylinder has charges induced on its ends.

wire, we could measure the speed of rotation of the system. We would have an “angular-velocity meter”!

But are you wondering: “What if I put myself in the frame of reference of the rotating cylinder? Then there is just a charged cylinder at rest, and I know that the electrostatic equations say there will be *no* electric fields inside, so there will be no force pushing charges to the center. So something must be wrong.” But there is nothing wrong. There is no “relativity of rotation.” A rotating system is *not* an inertial frame, and the laws of physics are different. We must be sure to use equations of electromagnetism only with respect to inertial coordinate systems.

It would be nice if we could measure the absolute rotation of the earth with such a charged cylinder, but unfortunately the effect is much too small to observe even with the most delicate instruments now available.

14-5 The field of a small loop; the magnetic dipole

Let’s use the vector-potential method to find the magnetic field of a small loop of current. As usual, by “small” we mean simply that we are interested in the fields only at distances large compared with the size of the loop. It will turn out that any small loop is a “magnetic dipole.” That is, it produces a *magnetic* field like the electric field from an electric dipole.

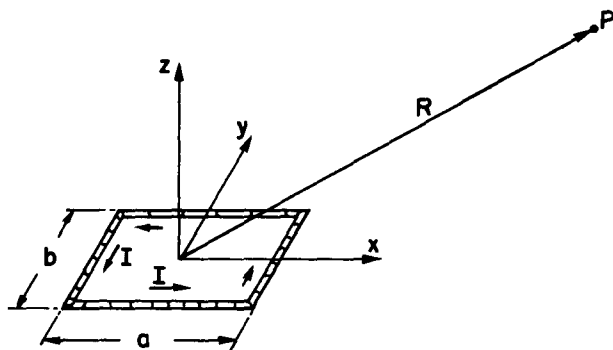


Fig. 14-6. A rectangular loop of wire with the current I . What is the magnetic field at P ? ($R \gg a$, or b .)

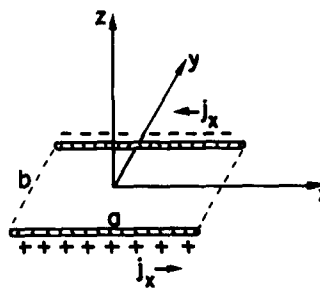


Fig. 14-7. The distribution of j_x in the current loop of Fig. 14-6.

We take first a rectangular loop, and choose our coordinates as shown in Fig. 14-6. There are no currents in the z -direction, so A_z is zero. There are currents in the x -direction on the two sides of length a . In each leg, the current density (and current) is uniform. So the solution for A_x is just like the electrostatic potential from two charged rods (see Fig. 14-7). Since the rods have opposite charges, their electric potential at large distances would be just the dipole potential (Section 6-5). At the point P in Fig. 14-6, the potential would be

$$\phi = \frac{1}{4\pi\epsilon_0} \frac{\mathbf{p} \cdot \mathbf{e}_R}{R^2}, \quad (14.28)$$

where \mathbf{p} is the dipole moment of the charge distribution. The dipole moment, in this case, is the total charge on one rod times the separation between them:

$$\mathbf{p} = \lambda ab. \quad (14.29)$$

The dipole moment points in the negative y -direction, so the cosine of the angle between \mathbf{R} and \mathbf{p} is $-y/R$ (where y is the coordinate of P). So we have

$$\phi = -\frac{1}{4\pi\epsilon_0} \frac{\lambda ab}{R^2} \frac{y}{R}.$$

We get A_x simply by replacing λ by I/c^2 :

$$A_x = -\frac{Iab}{4\pi\epsilon_0 c^2} \frac{y}{R^3}. \quad (14.30)$$

By the same reasoning,

$$A_y = \frac{Iab}{4\pi\epsilon_0 c^2} \frac{x}{R^3}. \quad (14.31)$$

Again, A_y is proportional to x and A_x is proportional to $-y$, so the vector potential (at large distances) goes in circles around the z -axis, circulating in the same sense as I in the loop, as shown in Fig. 14-8.

The strength of \mathcal{A} is proportional to Iab , which is the current times the area of the loop. This product is called the *magnetic dipole moment* (or, often, just “magnetic moment”) of the loop. We represent it by μ :

$$\mu = Iab. \quad (14.32)$$

The vector potential of a small plane loop of *any* shape (circle, triangle, etc.) is also given by Eqs. (14.30) and (14.31) provided we replace Iab by

$$\mu = I \cdot (\text{area of loop}). \quad (14.33)$$

We leave the proof of this to you.

We can put our equation in vector form if we define the direction of the vector μ to be the normal to the plane of the loop, with a positive sense given by the right-hand rule (Fig. 14-8). Then we can write

$$\mathcal{A} = \frac{1}{4\pi\epsilon_0 c^2} \frac{\mu \times \mathbf{R}}{R^3} = \frac{1}{4\pi\epsilon_0 c^2} \frac{\mu \times \mathbf{e}_R}{R^2}. \quad (14.34)$$

We have still to find \mathbf{B} . Using (14.33) and (14.34), together with (14.4), we get

$$B_x = -\frac{\partial}{\partial z} \frac{\mu}{4\pi\epsilon_0 c^2} \frac{x}{R^3} = \dots \frac{3xz}{R^5} \quad (14.35)$$

(where by \dots we mean $\mu/4\pi\epsilon_0 c^2$),

$$\begin{aligned} B_y &= \frac{\partial}{\partial z} \left(-\dots \frac{y}{R^3} \right) = \dots \frac{3yz}{R^5}, \\ B_z &= \frac{\partial}{\partial x} \left(\dots \frac{x}{R^3} \right) - \frac{\partial}{\partial y} \left(-\dots \frac{y}{R^3} \right) \\ &= -\dots \left(\frac{1}{r^3} - \frac{3z^2}{r^5} \right). \end{aligned} \quad (14.36)$$

The components of the \mathbf{B} -field behave exactly like those of the \mathbf{E} -field for a dipole oriented along the z -axis. (See Eqs. (6.14) and (6.15); also Fig. 6-5.) That’s why we call the loop a magnetic dipole. The word “dipole” is slightly misleading when applied to a magnetic field because there are *no* magnetic “poles” that correspond to electric charges. The magnetic “dipole field” is not produced by two “charges,” but by an elementary current loop.

It is curious, though, that starting with completely different laws, $\nabla \cdot \mathbf{E} = \rho/\epsilon_0$ and $\nabla \times \mathbf{B} = \mathbf{j}/\epsilon_0 c^2$, we can end up with the same kind of a field. Why should that be? It is because the dipole fields appear only when we are far away from all charges or currents. So through most of the relevant space the equations for \mathbf{E} and \mathbf{B} are identical: both have zero divergence and zero curl. So they give the same solutions. However, the *sources* whose configuration we summarize by the dipole moments are physically quite different—in one case, it’s a circulating current; in the other, a pair of charges, one above and one below the plane of the loop for the corresponding field.

14-6 The vector potential of a circuit

We are often interested in the magnetic fields produced by circuits of wire in which the diameter of the wire is very small compared with the dimensions of the whole system. In such cases, we can simplify the equations for the magnetic field.

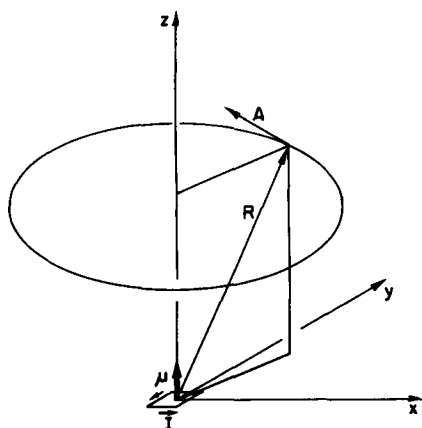


Fig. 14-8. The vector potential of a small current loop at the origin (in the xy -plane); a magnetic dipole field.

For a thin wire we can write our volume element as

$$dV = S ds,$$

where S is the cross-sectional area of the wire and ds is the element of distance along the wire. In fact, since the vector ds is in the same direction as \mathbf{j} , as shown in Fig. 14-9 (and we can assume that \mathbf{j} is constant across any given cross section), we can write a vector equation:

$$\mathbf{j} dV = jS ds. \quad (14.37)$$

But jS is just what we call the current I in a wire, so our integral for the vector potential (14.19) becomes

$$A(1) = \frac{1}{4\pi\epsilon_0 c^2} \int \frac{I ds_2}{r_{12}} \quad (14.38)$$

(see Fig. 14-10). (We assume that I is the same throughout the circuit. If there are several branches with different currents, we should, of course, use the appropriate I for each branch.)

Again, we can find the fields from (14.38) either by integrating directly or by solving the corresponding electrostatic problems.

14-7 The law of Biot and Savart

In studying electrostatics we found that the electric field of a known charge distribution could be obtained directly with an integral (Eq. 4-16):

$$\mathbf{E}(1) = \frac{1}{4\pi\epsilon_0} \int \frac{\rho(2)\mathbf{e}_{12} dV_2}{r_{12}^2}.$$

As we have seen, it is usually more work to evaluate this integral—there are really three integrals, one for each component—than to do the integral for the potential and take its gradient.

There is a similar integral which relates the magnetic field to the currents. We already have an integral for A , Eq. (14.19); we can get an integral for \mathbf{B} by taking the curl of both sides:

$$\mathbf{B}(1) = \nabla \times A(1) = \nabla \times \left[\frac{1}{4\pi\epsilon_0 c^2} \int \frac{j(2) dV_2}{r_{12}} \right]. \quad (14.39)$$

Now we must be careful: The curl operator means taking the derivatives of $A(1)$, that is, it operates only on the coordinates (x_1, y_1, z_1) . We can move the $\nabla \times$ operator inside the integral sign if we remember that it operates only on variables with the subscript 1, which of course, appear only in

$$r_{12} = [(x_1 - x_2)^2 + (y_1 - y_2)^2 + (z_1 - z_2)^2]^{1/2}. \quad (14.40)$$

We have, for the x -component of \mathbf{B} ,

$$\begin{aligned} B_x &= \frac{\partial A_z}{\partial y_1} - \frac{\partial A_y}{\partial z_1} \\ &= \frac{1}{4\pi\epsilon_0 c^2} \int \left[j_z \frac{\partial}{\partial y_1} \left(\frac{1}{r_{12}} \right) - j_y \frac{\partial}{\partial z_1} \left(\frac{1}{r_{12}} \right) \right] dV_2 \\ &= -\frac{1}{4\pi\epsilon_0 c^2} \int \left[j_z \frac{y_1 - y_2}{r_{12}^3} - j_y \frac{z_1 - z_2}{r_{12}^3} \right] dV_2. \end{aligned} \quad (14.41)$$

The quantity in brackets is just the x -component of

$$\frac{\mathbf{j} \times \mathbf{r}_{12}}{r_{12}^3} = \frac{\mathbf{j} \times \mathbf{e}_{12}}{r_{12}^2}.$$

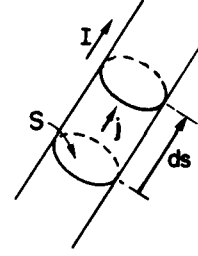


Fig. 14-9. For a fine wire $\mathbf{j} dV$ is the same as $I ds$.

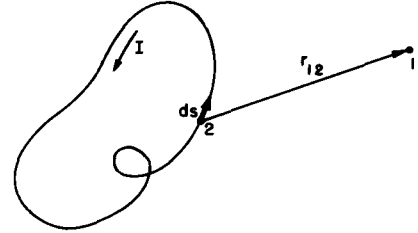


Fig. 14-10. The magnetic field of a wire can be obtained from an integral around the circuit.

Corresponding results will be found for the other components, so we have

$$\mathbf{B}(1) = \frac{1}{4\pi\epsilon_0 c^2} \int \frac{\mathbf{j}(2) \times \mathbf{e}_{12}}{r_{12}^2} dV_2. \quad (14.42)$$

The integral gives \mathbf{B} directly in terms of the known currents. The geometry involved is the same as that shown in Fig. 14-2.

If the currents exist only in circuits of small wires we can, as in the last section, immediately do the integral across the wire, replacing $\mathbf{j} dV$ by $I d\mathbf{s}$, where $d\mathbf{s}$ is an element of length of the wire. Then, using the symbols in Fig. 14-10,

$$\mathbf{B}(1) = -\frac{1}{4\pi\epsilon_0 c^2} \int \frac{I \mathbf{e}_{12} \times d\mathbf{s}_2}{r_{12}^2}. \quad (14.43)$$

(The minus sign appears because we have reversed the order of the cross product.) This equation for \mathbf{B} is called the *Biot-Savart law*, after its discoverers. It gives a formula for obtaining directly the magnetic field produced by wires carrying currents.

You may wonder: “What is the advantage of the vector potential if we can find \mathbf{B} directly with a vector integral? After all, \mathbf{A} also involves three integrals!” Because of the cross product, the integrals for \mathbf{B} are usually more complicated, as is evident from Eq. (14.41). Also, since the integrals for \mathbf{A} are like those of electrostatics, we may already know them. Finally, we will see that in more advanced theoretical matters (in relativity, in advanced formulations of the laws of mechanics, like the principle of least action to be discussed later, and in quantum mechanics) the vector potential plays an important role.

The Vector Potential

15-1 The forces on a current loop; energy of a dipole

In the last chapter we studied the magnetic field produced by a small rectangular current loop. We found that it is a dipole field, with the dipole moment given by

$$\mu = IA, \quad (15.1)$$

where I is the current and A is the area of the loop. The direction of the moment is normal to the plane of the loop, so we can also write

$$\mu = IAn,$$

where \mathbf{n} is the unit normal to the area A .

A current loop—or magnetic dipole—not only produces magnetic fields, but will also experience forces when placed in the magnetic field of other currents. We will look first at the forces on a rectangular loop in a uniform magnetic field. Let the z -axis be along the direction of the field, and the plane of the loop be placed through the y -axis, making the angle θ with the xy -plane as in Fig. 15-1. Then the magnetic moment of the loop—which is normal to its plane—will make the angle θ with the magnetic field.

Since the currents are opposite on opposite sides of the loop, the forces are also opposite, so there is no net force on the loop (when the field is uniform). Because of forces on the two sides marked 1 and 2 in the figure, however, there is a torque which tends to rotate the loop about the y -axis. The magnitude of these forces F_1 and F_2 is

$$F_1 = F_2 = IBb.$$

Their moment arm is

$$a \sin \theta,$$

so the torque is

$$\tau = Iab B \sin \theta,$$

or, since Iab is the magnetic moment of the loop,

$$\tau = \mu B \sin \theta.$$

The torque can be written in vector notation:

$$\tau = \mu \times B. \quad (15.2)$$

Although we have only shown that the torque is given by Eq. (15.2) in one rather special case, the result is right for a small loop of any shape, as we will see. You will remember that we found the same kind of relation for the torque on an electric dipole:

$$\tau = p \times E.$$

We now ask about the mechanical energy of our current loop. Since there is a torque, the energy evidently depends on the orientation. The principle of virtual work says that the torque is the rate of change of energy with angle, so we can write

$$dU = -\tau d\theta.$$

15-1 The forces on a current loop; energy of a dipole

15-2 Mechanical and electrical energies

15-3 The energy of steady currents

15-4 B versus A

15-5 The vector potential and quantum mechanics

15-6 What is true for statics is false for dynamics

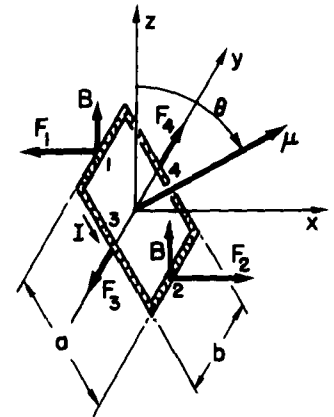


Fig. 15-1. A rectangular loop carrying the current I sits in a uniform field B (in the z -direction). The torque on the loop is $\tau = \mu \times B$, where the magnetic moment $\mu = Iab$.

Setting $\tau = -\mu B \sin \theta$, and integrating, we can write for the energy

$$U = -\mu B \cos \theta + \text{a constant.} \quad (15.3)$$

(The sign is negative because the torque tries to line up the moment with the field; the energy is lowest when μ and B are parallel.)

For reasons which we will discuss later, this energy is *not* the total energy of a current loop. (We have, for one thing, not taken into account the energy required to maintain the current in the loop.) We will, therefore, call this energy U_{mech} , to remind us that it is only part of the energy. Also, since we are leaving out some of the energy anyway, we can set the constant of integration equal to zero in Eq. (15.3). So we rewrite the equation:

$$U_{\text{mech}} = -\mu \cdot B. \quad (15.4)$$

Again, this corresponds to our result for an electric dipole:

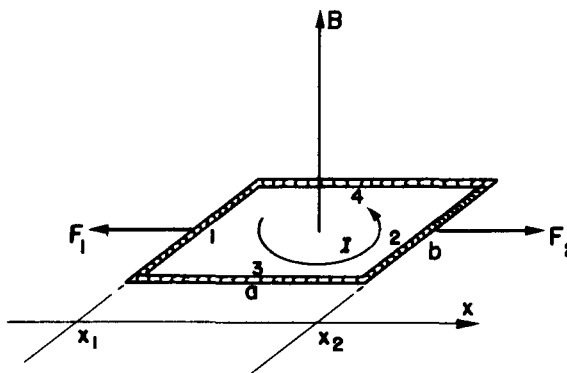
$$U = -p \cdot E. \quad (15.5)$$

Now the electrostatic energy U in Eq. (15.5) is the true energy, but U_{mech} in (15.4) is not the real energy. It *can*, however, be used in computing forces, by the principle of virtual work, supposing that the current in the loop—or at least μ —is kept constant.

We can show for our rectangular loop that U_{mech} also corresponds to the mechanical work done in bringing the loop into the field. The total force on the loop is zero only in a uniform field; in a nonuniform field there *are* net forces on a current loop. In putting the loop into a region with a field, we must have gone through places where the field was not uniform, and so work was done. To make the calculation simple, we shall imagine that the loop is brought into the field with its moment pointing along the field. (It can be rotated to its final position after it is in place.)

Imagine that we want to move the loop in the x -direction—toward a region of stronger field—and that the loop is oriented as shown in Fig. 15-2. We start somewhere where the field is zero and integrate the force times the distance as we bring the loop into the field.

Fig. 15-2. A loop is carried along the x -direction through the field B , at right angles to x .



First, let's compute the work done on each side separately and then take the sum (rather than adding the forces before integrating). The forces on sides 3 and 4 are at right angles to the direction of motion, so no work is done on them. The force on side 2 is $IbB(x)$ in the x -direction, and to get the work done against the magnetic forces we must integrate this from some x where the field is zero, say at $x = -\infty$, to x_2 , its present position:

$$W_2 = - \int_{-\infty}^{x_2} F_2 dx = -Ib \int_{-\infty}^{x_2} B(x) dx. \quad (15.6)$$

Similarly, the work done against the forces on side 1 is

$$W_1 = - \int_{-\infty}^{x_1} F_1 dx = Ib \int_{-\infty}^{x_1} B(x) dx. \quad (15.7)$$

To find each integral, we need to know how $B(x)$ depends on x . But notice that side 1 follows along right behind side 2, so that its integral includes most of the work done on side 2. In fact, the sum of (15.6) and (15.7) is just

$$W = -Ib \int_{x_1}^{x_2} B(x) dx. \quad (15.8)$$

But if we are in a region where B is nearly the same on both sides 1 and 2, we can write the integral as

$$\int_{x_1}^{x_2} B(x) dx = (x_2 - x_1)B = aB,$$

where B is the field at the center of the loop. The total mechanical energy we have put in is

$$U_{\text{mech}} = W = -Iab B = -\mu B. \quad (15.9)$$

The result agrees with the energy we took for Eq. (15.4).

We would, of course, have gotten the same result if we had added the forces on the loop before integrating to find the work. If we let B_1 be the field at side 1 and B_2 be the field at side 2, then the total force in the x -direction is

$$F_x = Ib(B_2 - B_1).$$

If the loop is “small,” that is, if B_2 and B_1 are not too different, we can write

$$B_2 = B_1 + \frac{\partial B}{\partial x} \Delta x = B_1 + \frac{\partial B}{\partial x} a.$$

So the force is

$$F_x = Iab \frac{\partial B}{\partial x}. \quad (15.10)$$

The total work done on the loop by *external* forces is

$$-\int_{-\infty}^{\infty} F_x dx = -Iab \int \frac{\partial B}{\partial x} dx = -IabB,$$

which is again just $-\mu B$. Only now we see why it is that the *force* on a small current loop is proportional to the derivative of the magnetic field, as we would expect from

$$F_x \Delta x = -\Delta U_{\text{mech}} = -\Delta(-\mu \cdot \mathbf{B}). \quad (15.11)$$

Our result, then, is that even though $U_{\text{mech}} = -\mu \cdot \mathbf{B}$ may not include all the energy of a system—it is a fake kind of energy—it can still be used with the principle of virtual work to find the forces on steady current loops.

15-2 Mechanical and electrical energies

We want now to show why the energy U_{mech} discussed in the previous section is not the correct energy associated with steady currents—that it does not keep track of the total energy in the world. We have, indeed, emphasized that it can be used like the energy, for computing forces from the principle of virtual work, *provided* that the current in the loop (and all *other* currents) do not change. Let’s see why all this works.

Imagine that the loop in Fig. 15-2 is moving in the $+x$ -direction and take the z -axis in the direction of \mathbf{B} . The conduction electrons in side 2 will experience a force along the wire, in the y -direction. But because of their flow—as an electric current—there is a component of their motion in the same direction as the force. Each electron is, therefore, having work done on it at the rate $F_y v_y$, where v_y is the component of the electron velocity along the wire. We will call this work done on the electrons *electrical* work. Now it turns out that if the loop is moving in a *uniform* field, the total electrical work is zero, since positive work is done on some parts of the loop and an equal amount of negative work is done on other parts.

But this is not true if the circuit is moving in a nonuniform field—then there *will* be a net amount of work done on the electrons. In general, this work would tend to change the flow of the electrons, but if the current is being held constant, energy must be absorbed or delivered by the battery or other source that is keeping the current steady. This energy was not included when we computed U_{mech} in Eq. (15.9), because our computations included only the mechanical forces on the body of the wire.

You may be thinking: But the force on the electrons depends on how *fast* the wire is moved; perhaps if the wire is moved slowly enough this electrical energy can be neglected. It is true that the *rate* at which the electrical energy is delivered is proportional to the speed of the wire, but the *total* energy delivered is proportional also to the *time* that this rate goes on. So the total electrical energy is proportional to the velocity times the time, which is just the distance moved. For a given distance moved in a field the same amount of electrical work is done.

Let's consider a segment of wire of unit length carrying the current I and moving in a direction perpendicular to itself and to a magnetic field \mathbf{B} with the speed v_{wire} . Because of the current the electrons will have a drift velocity v_{drift} along the wire. The component of the magnetic force on each electron in the direction of the drift is $q_e v_{\text{wire}} B$. So the rate at which electrical work is being done is $F v_{\text{drift}} = (q_e v_{\text{wire}} B) v_{\text{drift}}$. If there are N conduction electrons in the unit length of the wire, the total rate at which electrical work is being done is

$$\frac{dU_{\text{elect}}}{dt} = N q_e v_{\text{wire}} B v_{\text{drift}}.$$

But $N q_e v_{\text{drift}} = I$, the current in the wire, so

$$\frac{dU_{\text{elect}}}{dt} = I v_{\text{wire}} B.$$

Now since the current is held constant, the forces on the conduction electrons do not cause them to accelerate; the electrical energy is not going into the electrons but into the source that is keeping the current constant.

But notice that the force on the *wire* is IB , so $IB v_{\text{wire}}$ is also the rate of *mechanical work* done on the wire, $dU_{\text{mech}}/dt = IB v_{\text{wire}}$. We conclude that the mechanical work done on the wire is just equal to the electrical work done on the current source, so the energy of the loop is *a constant!*

This is not a coincidence, but a consequence of the law we already know. The total force on each charge in the wire is

$$\mathbf{F} = q(\mathbf{E} + \mathbf{v} \times \mathbf{B}).$$

The rate at which work is done is

$$\mathbf{v} \cdot \mathbf{F} = q[\mathbf{v} \cdot \mathbf{E} + \mathbf{v} \cdot (\mathbf{v} \times \mathbf{B})]. \quad (15.12)$$

If there are no electric fields we have only the second term, which is always zero. We shall see later that *changing* magnetic fields produce electric fields, so our reasoning applies only to moving wires in steady magnetic fields.

How is it then that the principle of virtual work gives the right answer? Because we *still* have not taken into account the *total* energy of the world. We have not included the energy of the currents that are *producing* the magnetic field we start out with.

Suppose we imagine a complete system such as that drawn in Fig. 15-3(a), in which we are moving our loop with the current I_1 into the magnetic field \mathbf{B}_1 produced by the current I_2 in a coil. Now the current I_1 in the loop will also be producing some magnetic field \mathbf{B}_2 at the coil. If the loop is moving, the field \mathbf{B}_2 will be changing. As we shall see in the next chapter, a changing magnetic field generates an \mathbf{E} -field; and this \mathbf{E} -field will do work on the charges in the coil. This energy must also be included in our balance sheet of the total energy.

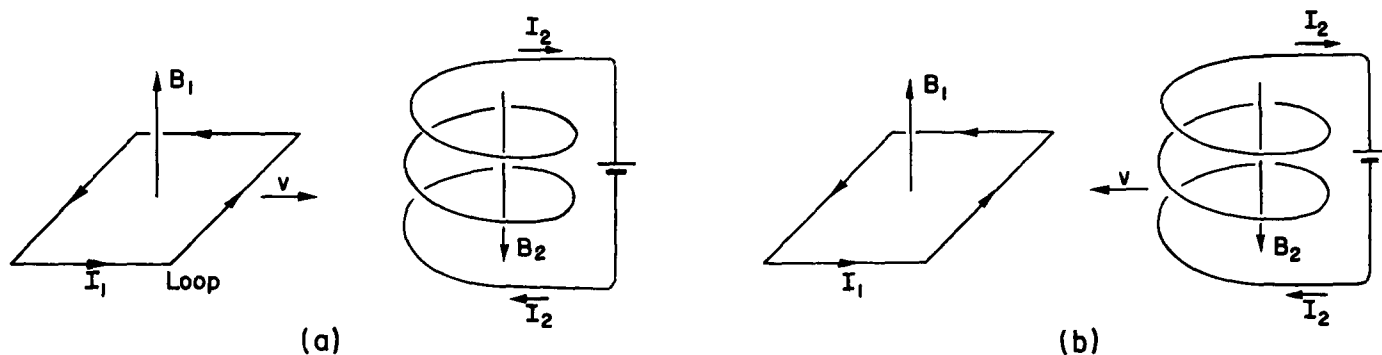


Fig. 15-3. Finding the energy of a small loop in a magnetic field.

We could wait until the next chapter to find out about this new energy term, but we can also see what it will be if we use the principle of relativity in the following way. When we are moving the loop toward the stationary coil we know that its electrical energy is just equal and opposite to the mechanical work done. So

$$U_{\text{mech}} + U_{\text{elect}}(\text{loop}) = 0.$$

Suppose now we look at what is happening from a different point of view, in which the loop is at rest, and the coil is moved toward it. The coil is then moving into the field produced by the loop. The same arguments would give that

$$U_{\text{mech}} + U_{\text{elect}}(\text{coil}) = 0.$$

The mechanical energy is the same in the two cases because it comes from the force between the two circuits.

The sum of the two equations gives

$$2U_{\text{mech}} + U_{\text{elect}}(\text{loop}) + U_{\text{elect}}(\text{coil}) = 0.$$

The total energy of the whole system is, of course, the sum of the two electrical energies plus the mechanical energy taken only *once*. So we have

$$U_{\text{total}} = U_{\text{elect}}(\text{loop}) + U_{\text{elect}}(\text{coil}) + U_{\text{mech}} = -U_{\text{mech}}. \quad (15.13)$$

The total energy of the world is really the *negative* of U_{mech} . If we want the true energy of a magnetic dipole, for example, we should write

$$U_{\text{total}} = +\boldsymbol{\mu} \cdot \mathbf{B}.$$

It is only if we make the condition that all currents are constant that we can use only a part of the energy, U_{mech} (which is always the negative of the true energy), to find the mechanical forces. In a more general problem, we must be careful to include all energies.

We have seen an analogous situation in electrostatics. We showed that the energy of a capacitor is equal to $Q^2/2C$. When we use the principle of virtual work to find the force between the plates of the capacitor, the change in energy is equal to $Q^2/2$ times the change in $1/C$. That is,

$$\Delta U = \frac{Q^2}{2} \Delta \left(\frac{1}{C} \right) = -\frac{Q^2}{2} \frac{\Delta C}{C^2}. \quad (15.14)$$

Now suppose that we were to calculate the work done in moving two conductors subject to the different condition that the voltage between them is held constant. Then we can get the right answers for force from the principle of virtual work if we do something artificial. Since $Q = CV$, the real energy is $\frac{1}{2}CV^2$. But if we define an artificial energy equal to $-\frac{1}{2}CV^2$, then the principle of virtual work can be used to get forces by setting the change in the artificial energy equal to the

mechanical work, provided that we insist that the voltage V be held constant. Then

$$\Delta U_{\text{mech}} = \Delta \left(-\frac{CV^2}{2} \right) = -\frac{V^2}{2} \Delta C, \quad (15.15)$$

which is the same as Eq. (15.14). We get the correct result even though we are neglecting the work done by the electrical system to keep the voltage constant. Again, this electrical energy is just twice as big as the mechanical energy and of the opposite sign.

Thus if we calculate artificially, disregarding the fact that the source of the potential has to do work to maintain the voltages constant, we get the right answer. It is exactly analogous to the situation in magnetostatics.

15-3 The energy of steady currents

We can now use our knowledge that $U_{\text{total}} = -U_{\text{mech}}$ to find the true energy of steady currents in magnetic fields. We can begin with the true energy of a small current loop. Calling U_{total} just U , we write

$$U = \mu \cdot B. \quad (15.16)$$

Although we calculated this energy for a plane rectangular loop, the same result holds for a small plane loop of any shape.

We can find the energy of a circuit of any shape by imagining that it is made up of small current loops. Say we have a wire in the shape of the loop Γ of Fig. 15-4. We fill in this curve with the surface S , and on the surface mark out a large number of small loops, each of which can be considered plane. If we let the current I circulate around *each* of the little loops, the net result will be the same as a current around Γ , since the currents will cancel on all lines internal to Γ . Physically, the system of little currents is indistinguishable from the original circuit. The energy must also be the same, and so is just the sum of the energies of the little loops.

If the area of each little loop is Δa , its energy is $I \Delta a B_n$, where B_n is the component normal to Δa . The total energy is

$$U = \sum I B_n \Delta a.$$

Going to the limit of infinitesimal loops, the sum becomes an integral, and

$$U = I \int B_n da = I \int \mathbf{B} \cdot \mathbf{n} da, \quad (15.17)$$

where \mathbf{n} is the unit normal to da .

If we set $\mathbf{B} = \nabla \times \mathbf{A}$, we can connect the surface integral to a line integral, using Stokes' theorem,

$$I \int_S (\nabla \times \mathbf{A}) \cdot \mathbf{n} da = I \oint_{\Gamma} \mathbf{A} \cdot d\mathbf{s}, \quad (15.18)$$

where $d\mathbf{s}$ is the line element along Γ . So we have the energy for a circuit of any shape:

$$U = I \oint_{\text{circuit}} \mathbf{A} \cdot d\mathbf{s}. \quad (15.19)$$

In this expression \mathbf{A} refers, of course, to the vector potential due to those currents (other than the I in the wire) which produce the field \mathbf{B} at the wire.

Now any distribution of steady currents can be imagined to be made up of filaments that run parallel to the lines of current flow. For each pair of such circuits, the energy is given by (15.19), where the integral is taken around one circuit, using the vector potential \mathbf{A} from the other circuit. For the total energy we want the sum of all such pairs. If, instead of keeping track of the pairs, we take the complete sum over all the filaments, we would be counting the energy twice (we saw a similar effect in electrostatics), so the total energy can be written

$$U = \frac{1}{2} \int \mathbf{j} \cdot \mathbf{A} dV. \quad (15.20)$$

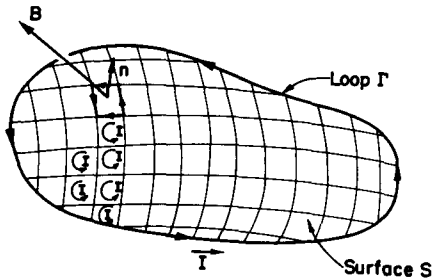


Fig. 15-4. The energy of a large loop in a magnetic field can be considered as the sum of energies of smaller loops.

This formula corresponds to the result we found for the electrostatic energy:

$$U = \frac{1}{2} \int \rho \phi dV. \quad (15.21)$$

So we may if we wish think of A as a kind of potential energy for currents in magnetostatics. Unfortunately, this idea is not too useful, because it is true only for static fields. In fact, neither of the equations (15.20) and (15.21) gives the correct energy when the fields change with time.

15-4 B versus A

In this section we would like to discuss the following questions: Is the vector potential merely a device which is useful in making calculations—as the scalar potential is useful in electrostatics—or is the vector potential a “real” field? Isn’t the magnetic field the “real” field, because it is responsible for the force on a moving particle? First we should say that the phrase “a real field” is not very meaningful. For one thing, you probably don’t feel that the magnetic field is very “real” anyway, because even the whole idea of a field is a rather abstract thing. You cannot put out your hand and feel the magnetic field. Furthermore, the value of the magnetic field is not very definite; by choosing a suitable moving coordinate system, for instance, you can make a magnetic field at a given point disappear.

What we mean here by a “real” field is this: a real field is a mathematical function we use for avoiding the idea of action at a distance. If we have a charged particle at the position P , it is affected by other charges located at some distance from P . One way to describe the interaction is to say that the other charges make some “condition”—whatever it may be—in the environment at P . If we know that condition, which we describe by giving the electric and magnetic fields, then we can determine completely the behavior of the particle—with no further reference to how those conditions came about.

In other words, if those other charges were altered in some way, but the conditions at P that are described by the electric and magnetic field at P remain the same, then the motion of the charge will also be the same. A “real” field is then a set of numbers we specify in such a way that what happens *at a point* depends only on the numbers *at that point*. We do not need to know any more about what’s going on at other places. It is in this sense that we will discuss whether the vector potential is a “real” field.

You may be wondering about the fact that the vector potential is not unique—that it can be changed by adding the gradient of any scalar with no change at all in the forces on particles. That has not, however, anything to do with the question of reality in the sense that we are talking about. For instance, the magnetic field is in a sense altered by a relativity change (as are also E and A). But we are not worried about what happens if the field *can* be changed in this way. That doesn’t really make any difference; that has nothing to do with the question of whether the vector potential is a proper “real” field for describing magnetic effects, or whether it is just a useful mathematical tool.

We should also make some remarks on the usefulness of the vector potential A . We have seen that it can be used in a formal procedure for calculating the magnetic fields of known currents, just as ϕ can be used to find electric fields. In electrostatics we saw that ϕ was given by the scalar integral

$$\phi(1) = \frac{1}{4\pi\epsilon_0} \int \frac{\rho(2)}{r_{12}} dV_2. \quad (15.22)$$

From this ϕ , we get the three components of E by three differential operations. This procedure is usually easier to handle than evaluating the three integrals in the vector formula

$$E(1) = \frac{1}{4\pi\epsilon_0} \int \frac{\rho(2)e_{12}}{r_{12}^2} dV_2. \quad (15.23)$$

First, there are three integrals; and second, each integral is in general somewhat more difficult.

The advantages are much less clear for magnetostatics. The integral for \mathcal{A} is already a vector integral:

$$\mathcal{A}(1) = \frac{1}{4\pi\epsilon_0 c^2} \int \frac{\mathbf{j}(2) dV_2}{r_{12}}, \quad (15.24)$$

which is, of course, three integrals. Also, when we take the curl of \mathcal{A} to get \mathbf{B} , we have six derivatives to do and combine by pairs. It is not immediately obvious whether in most problems this procedure is really any easier than computing \mathbf{B} directly from

$$\mathbf{B}(1) = \frac{1}{4\pi\epsilon_0 c^2} \int \frac{\mathbf{j}(2) \times \mathbf{e}_{12}}{r_{12}^2} dV_2. \quad (15.25)$$

Using the vector potential is often more difficult for simple problems for the following reason. Suppose we are interested only in the magnetic field \mathbf{B} at one point, and that the problem has some nice symmetry—say we want the field at a point on the axis of a ring of current. Because of the symmetry, we can easily get \mathbf{B} by doing the integral of Eq. (15.25). If, however, we were to find \mathcal{A} first, we would have to compute \mathbf{B} from *derivatives* of \mathcal{A} , so we must know what \mathcal{A} is at all points in the *neighborhood* of the point of interest. And most of these points are off the axis of symmetry, so the integral for \mathcal{A} gets complicated. In the ring problem, for example, we would need to use elliptic integrals. In such problems, \mathcal{A} is clearly not very useful. It is true that in many complex problems it is easier to work with \mathcal{A} , but it would be hard to argue that this ease of technique would justify making you learn about one more vector field.

We have introduced \mathcal{A} because it *does* have an important physical significance. Not only is it related to the energies of currents, as we saw in the last section, but it is also a “real” physical field in the sense that we described above. In classical mechanics it is clear that we can write the force on a particle as

$$\mathbf{F} = q(\mathbf{E} + \mathbf{v} \times \mathbf{B}), \quad (15.26)$$

so that, given the forces, everything about the motion is determined. In any region where $\mathbf{B} = 0$ even if \mathcal{A} is not zero, such as outside a solenoid, there is no discernible effect of \mathcal{A} . Therefore for a long time it was believed that \mathcal{A} was not a “real” field. It turns out, however, that there are phenomena involving quantum mechanics which show that the field \mathcal{A} is in fact a “real” field in the sense we have defined it. In the next section we will show you how that works.

15-5 The vector potential and quantum mechanics

There are many changes in what concepts are important when we go from classical to quantum mechanics. We have already discussed some of them in Vol. I. In particular, the force concept gradually fades away, while the concepts of energy and momentum become of paramount importance. You remember that instead of particle motions, one deals with probability amplitudes which vary in space and time. In these amplitudes there are wavelengths related to momenta, and frequencies related to energies. The momenta and energies, which determine the phases of wave functions, are therefore the important quantities in quantum mechanics. Instead of forces, we deal with the way interactions change the wavelength of the waves. The idea of a force becomes quite secondary—if it is there at all. When people talk about nuclear forces, for example, what they usually analyze and work with are the energies of interaction of two nucleons, and not the force between them. Nobody ever differentiates the energy to find out what the force looks like. In this section we want to describe how the vector and scalar potentials enter into quantum mechanics. It is, in fact, just because momentum and energy play a central role in quantum mechanics that \mathcal{A} and ϕ provide the most direct way of introducing electromagnetic effects into quantum descriptions.

We must review a little how quantum mechanics works. We will consider again the imaginary experiment described in Chapter 37 of Vol. I, in which elec-

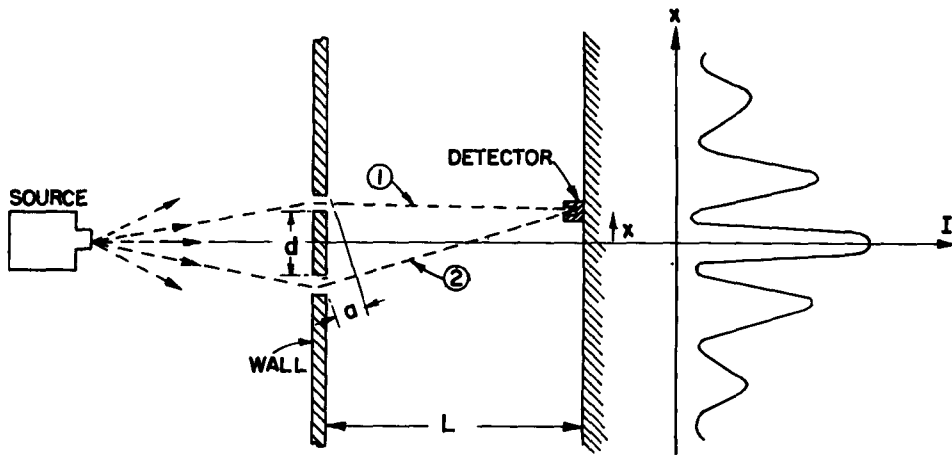


Fig. 15-5. An interference experiment with electrons (see also Chapter 37 of Vol. I).

trons are diffracted by two slits. The arrangement is shown again in Fig. 15-5. Electrons, all of nearly the same energy, leave the source and travel toward a wall with two narrow slits. Beyond the wall is a "backstop" with a movable detector. The detector measures the rate, which we call I , at which electrons arrive at a small region of the backstop at the distance x from the axis of symmetry. The rate is proportional to the probability that an individual electron that leaves the source will reach that region of the backstop. This probability has the complicated-looking distribution shown in the figure, which we understand as due to the interference of two amplitudes, one from each slit. The interference of the two amplitudes depends on their phase difference. That is, if the amplitudes are $C_1 e^{i\Phi_1}$ and $C_2 e^{i\Phi_2}$, the phase difference $\delta = \Phi_1 - \Phi_2$ determines their interference pattern [see Eq. (29.12) in Vol. I]. If the distance between the screen and the slits is L , and if the difference in the path lengths for electrons going through the two slits is a , as shown in the figure, then the phase difference of the two waves is given by

$$\delta = \frac{a}{\lambda}. \quad (15.27)$$

As usual, we let $\lambda = \lambda/2\pi$, where λ is the wavelength of the space variation of the probability amplitude. For simplicity, we will consider only values of x much less than L ; then we can set

$$a = \frac{x}{L} d$$

and

$$\delta = \frac{x}{L} \frac{d}{\lambda}. \quad (15.28)$$

When x is zero, δ is zero; the waves are in phase, and the probability has a maximum. When δ is π , the waves are out of phase, they interfere destructively, and the probability is a minimum. So we get the wavy function for the electron intensity.

Now we would like to state the law that for quantum mechanics replaces the force law $\mathbf{F} = q\mathbf{v} \times \mathbf{B}$. It will be the law that determines the behavior of quantum-mechanical particles in an electromagnetic field. Since what happens is determined by amplitudes, the law must tell us how the magnetic influences affect the amplitudes; we are no longer dealing with the acceleration of a particle. The law is the following: the phase of the amplitude to arrive via any trajectory is changed by the presence of a magnetic field by an amount equal to the integral of the vector potential along the whole trajectory times the charge of the particle over Planck's constant. That is,

$$\text{Magnetic change in phase} = \frac{q}{\hbar} \int_{\text{trajectory}} \mathbf{A} \cdot d\mathbf{s}. \quad (15.29)$$

If there were no magnetic field there would be a certain phase of arrival. If there is a magnetic field anywhere, the phase of the arriving wave is increased by the integral in Eq. (15.29).

Although we will not need to use it for our present discussion, we mention that the effect of an electrostatic field is to produce a phase change given by the *negative* of the *time* integral of the scalar potential ϕ :

$$\text{Electric change in phase} = -\frac{q}{\hbar} \int \phi \, dt.$$

These two expressions are correct not only for static fields, but together give the correct result for *any* electromagnetic field, static or dynamic. This is the law that replaces $\mathbf{F} = q(\mathbf{E} + \mathbf{v} \times \mathbf{B})$. We want now, however, to consider only a static magnetic field.

Suppose that there is a magnetic field present in the two-slit experiment. We want to ask for the phase of arrival at the screen of the two waves whose paths pass through the two slits. Their interference determines where the maxima in the probability will be. We may call Φ_1 the phase of the wave along trajectory (1). If $\Phi_1(B = 0)$ is the phase without the magnetic field, then when the field is turned on the phase will be

$$\Phi_1 = \Phi_1(B = 0) + \frac{q}{\hbar} \int_{(1)} \mathbf{A} \cdot d\mathbf{s}. \quad (15.30)$$

Similarly, the phase for trajectory (2) is

$$\Phi_2 = \Phi_2(B = 0) + \frac{q}{\hbar} \int_{(2)} \mathbf{A} \cdot d\mathbf{s}. \quad (15.31)$$

The interference of the waves at the detector depends on the phase difference

$$\delta = \Phi_1(B = 0) - \Phi_2(B = 0) + \frac{q}{\hbar} \int_{(1)} \mathbf{A} \cdot d\mathbf{s} - \frac{q}{\hbar} \int_{(2)} \mathbf{A} \cdot d\mathbf{s}. \quad (15.32)$$

The no-field difference we will call $\delta(B = 0)$; it is just the phase difference we have calculated above in Eq. (15.28). Also, we notice that the two integrals can be written as *one* integral that goes forward along (1) and back along (2); we call this the closed path (1-2). So we have

$$\delta = \delta(B = 0) + \frac{q}{\hbar} \oint_{(1-2)} \mathbf{A} \cdot d\mathbf{s}. \quad (15.33)$$

This equation tells us how the electron motion is changed by the magnetic field; with it we can find the new positions of the intensity maxima and minima at the backstop.

Before we do that, however, we want to raise the following interesting and important point. You remember that the vector potential function has some arbitrariness. Two different vector potential functions \mathbf{A} and \mathbf{A}' whose difference is the gradient of some scalar function $\nabla\psi$, both represent the same magnetic field, since the curl of a gradient is zero. They give, therefore, the same classical force $q\mathbf{v} \times \mathbf{B}$. If in quantum mechanics the effects depend on the vector potential, *which* of the many possible \mathbf{A} -functions is correct?

The answer is that the same arbitrariness in \mathbf{A} continues to exist for quantum mechanics. If in Eq. (15.33) we change \mathbf{A} to $\mathbf{A}' = \mathbf{A} + \nabla\psi$, the integral on \mathbf{A} becomes

$$\oint_{(1-2)} \mathbf{A}' \cdot d\mathbf{s} = \oint_{(1-2)} \mathbf{A} \cdot d\mathbf{s} + \oint_{(1-2)} \nabla\psi \cdot d\mathbf{s}.$$

The integral of $\nabla\psi$ is around the *closed* path (1-2), but the integral of the tangential component of a gradient on a closed path is always zero, by Stokes' theorem. Therefore both \mathbf{A} and \mathbf{A}' give the same phase differences and the same quantum-mechanical interference effects. In both classical and quantum theory it is only the curl of \mathbf{A} that matters; any choice of the function of \mathbf{A} which has the correct curl gives the correct physics.

The same conclusion is evident if we use the results of Section 14-1. There we found that the line integral of \mathbf{A} around a closed path is the flux of \mathbf{B} through the path, which here is the flux between paths (1) and (2). Equation (15.33) can, if we wish, be written as

$$\delta = \delta(B = 0) + \frac{q}{\hbar} [\text{flux of } \mathbf{B} \text{ between (1) and (2)}], \quad (15.34)$$

where by the flux of \mathbf{B} we mean, as usual, the surface integral of the normal component of \mathbf{B} . The result depends only on \mathbf{B} , and therefore only on the curl of \mathbf{A} .

Now because we can write the result in terms of \mathbf{B} as well as in terms of \mathbf{A} , you might be inclined to think that the \mathbf{B} holds its own as a “real” field and that the \mathbf{A} can still be thought of as an artificial construction. But the definition of “real” field that we originally proposed was based on the idea that a “real” field would not act on a particle from a distance. We can, however, give an example in which \mathbf{B} is zero—or at least arbitrarily small—at any place where there is some chance to find the particles, so that it is not possible to think of it acting *directly* on them.

You remember that for a long solenoid carrying an electric current there is a \mathbf{B} -field inside but none outside, while there is lots of \mathbf{A} circulating around outside, as shown in Fig. 15-6. If we arrange a situation in which electrons are to be found only *outside* of the solenoid—only where there is \mathbf{A} —there will still be an influence on the motion, according to Eq. (15.33). Classically, that is impossible. Classically, the force depends only on \mathbf{B} ; in order to know that the solenoid is carrying current, the particle must go through it. But quantum-mechanically you can find out that there is a magnetic field inside the solenoid by going *around* it—without ever going close to it!

Suppose that we put a very long solenoid of small diameter just behind the wall and between the two slits, as shown in Fig. 15-7. The diameter of the solenoid is to be much smaller than the distance d between the two slits. In these circumstances, the diffraction of the electrons at the slit gives no appreciable probability that the electrons will get near the solenoid. What will be the effect on our interference experiment?

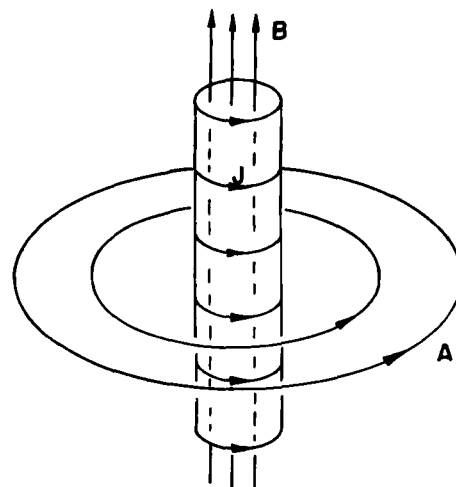


Fig. 15-6. The magnetic field and vector potential of a long solenoid.

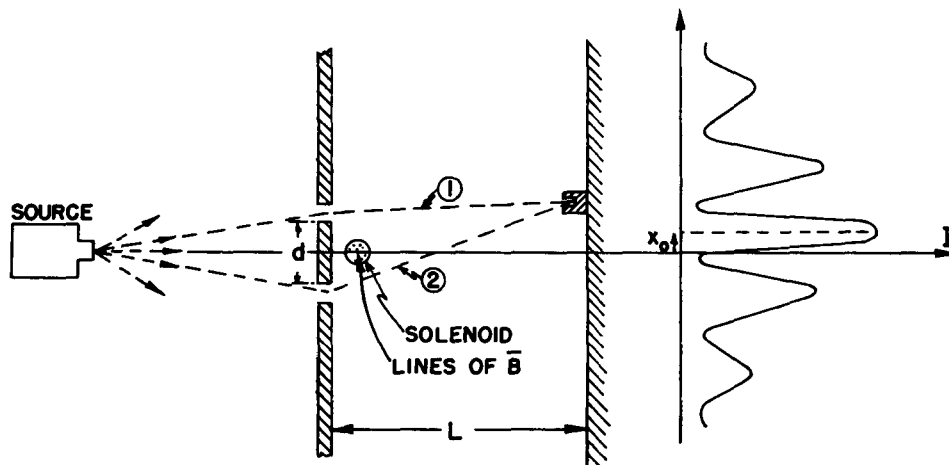


Fig. 15-7. A magnetic field can influence the motion of electrons even though it exists only in regions where there is an arbitrarily small probability of finding the electrons.

We compare the situation with and without a current through the solenoid. If we have no current, we have no \mathbf{B} or \mathbf{A} and we get the original pattern of electron intensity at the backstop. If we turn the current on in the solenoid and build up a magnetic field \mathbf{B} inside, then there is an \mathbf{A} outside. There is a shift in the phase difference proportional to the circulation of \mathbf{A} outside the solenoid, which will mean that the pattern of maxima and minima is shifted to a new position. In fact, since the flux of \mathbf{B} inside is a constant for any pair of paths, so also is the circulation of \mathbf{A} . For every arrival point there is the same phase change; this corresponds

to shifting the entire pattern in x by a constant amount, say x_0 , that we can easily calculate. The maximum intensity will occur where the phase difference between the two waves is zero. Using Eq. (15.32) or Eq. (15.33) for δ and Eq. (15.28) for $\delta(B = 0)$, we have

$$x_0 = -\frac{L}{d} \lambda \frac{q}{\hbar} \oint_{(1-2)} \mathbf{A} \cdot d\mathbf{s}, \quad (15.35)$$

or

$$x_0 = -\frac{L}{d} \lambda \frac{q}{\hbar} [\text{flux of } \mathbf{B} \text{ between (1) and (2)}]. \quad (15.36)$$

The pattern with the solenoid in place should appear* as shown in Fig. 15-7. At least, that is the prediction of quantum mechanics.

Precisely this experiment has recently been done. It is a very, very difficult experiment. Because the wavelength of the electrons is so small, the apparatus must be on a tiny scale to observe the interference. The slits must be very close together, and that means that one needs an exceedingly small solenoid. It turns out that in certain circumstances, iron crystals will grow in the form of very long, microscopically thin filaments called whiskers. When these iron whiskers are magnetized they are like a tiny solenoid, and there is no field outside except near the ends. The electron interference experiment was done with such a whisker between two slits, and the predicted displacement in the pattern of electrons was observed.

In our sense then, the \mathbf{A} -field is "real." You may say: "But there *was* a magnetic field." There was, but remember our original idea—that a field is "real" if it is what must be specified *at the position* of the particle in order to get the motion. The \mathbf{B} -field in the whisker acts at a distance. If we want to describe its influence not as action-at-a-distance, we must use the vector potential.

This subject has an interesting history. The theory we have described was known from the beginning of quantum mechanics in 1926. The fact that the vector potential appears in the wave equation of quantum mechanics (called the Schrödinger equation) was obvious from the day it was written. That it cannot be replaced by the magnetic field in any easy way was observed by one man after the other who tried to do so. This is also clear from our example of electrons moving in a region where there is no field and being affected nevertheless. But because in classical mechanics \mathbf{A} did not appear to have any direct importance and, furthermore, because it could be changed by adding a gradient, people repeatedly said that the vector potential had no direct physical significance—that only the magnetic and electric fields are "right" even in quantum mechanics. It seems strange in retrospect that no one thought of discussing this experiment until 1956, when Bohm and Aharonov first suggested it and made the whole question crystal clear. The implication was there all the time, but no one paid attention to it. Thus many people were rather shocked when the matter was brought up. That's why someone thought it would be worth while to do the experiment to see that it really was right, even though quantum mechanics, which had been believed for so many years, gave an unequivocal answer. It is interesting that something like this can be around for thirty years but, because of certain prejudices of what is and is not significant, continues to be ignored.

Now we wish to continue in our analysis a little further. We will show the connection between the quantum-mechanical formula and the classical formula—to show why it turns out that if we look at things on a large enough scale it will look as though the particles are acted on by a force equal to $q\mathbf{v} \times$ the curl of \mathbf{A} . To get classical mechanics from quantum mechanics, we need to consider cases in which all the wavelengths are very small compared with distances over which external conditions, like fields, vary appreciably. We shall not prove the result in great generality, but only in a very simple example, to show how it works. Again we consider the same slit experiment. But instead of putting all the magnetic field in a very tiny region between the slits, we imagine a magnetic field that extends

* If the field \mathbf{B} comes out of the plane of the figure, the flux as we have defined it is negative and x_0 is positive.

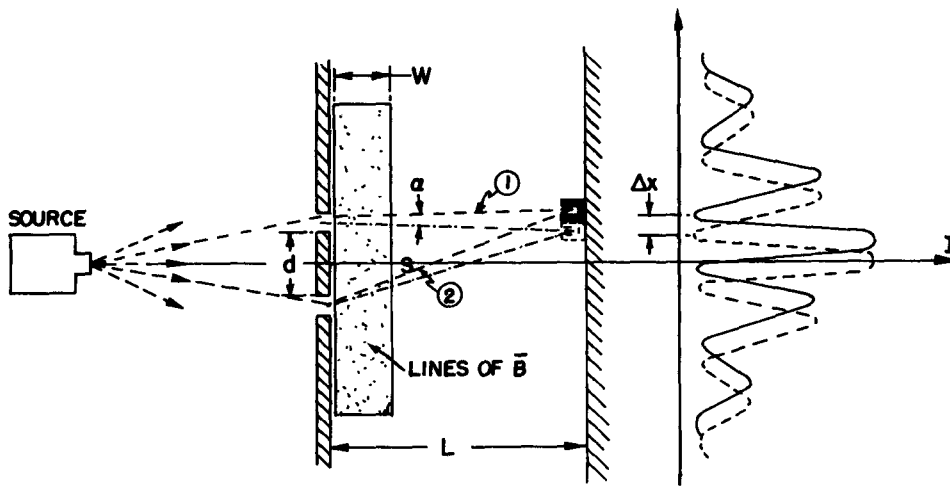


Fig. 15-8. The shift of the interference pattern due to a strip of magnetic field.

over a larger region behind the slits, as shown in Fig. 15-8. We will take the idealized case where we have a magnetic field which is uniform in a narrow strip of width w , considered small as compared with L . (That can easily be arranged; the backstop can be put as far out as we want.) In order to calculate the shift in phase, we must take the two integrals of A along the two trajectories (1) and (2). They differ, as we have seen, merely by the flux of B between the paths. To our approximation, the flux is Bwd . The phase difference for the two paths is then

$$\delta = \delta(B = 0) + \frac{q}{\hbar} Bwd. \quad (15.37)$$

We note that, to our approximation, the phase shift is independent of the angle. So again the effect will be to shift the whole pattern upward by an amount Δx . Using Eq. (15.28),

$$\Delta x = \frac{L\lambda}{d} \Delta\delta = \frac{L\lambda}{d} [\delta - \delta(B = 0)].$$

Using (15.37) for $\delta - \delta(B = 0)$,

$$\Delta x = L\lambda \frac{q}{\hbar} Bw. \quad (15.38)$$

Such a shift is equivalent to deflecting all the trajectories by the small angle α (see Fig. 15-8), where

$$\alpha = \frac{\Delta x}{L} = \frac{\lambda}{\hbar} qBw. \quad (15.39)$$

Now classically we would also expect a thin strip of magnetic field to deflect all trajectories through some small angle, say α' , as shown in Fig. 15-9(a). As the electrons go through the magnetic field, they feel a transverse force $qv \times B$ which lasts for a time w/v . The change in their transverse momentum is just equal to this impulse, so

$$\Delta p_x = qwB. \quad (15.40)$$

The angular deflection [Fig. 15-9(b)] is equal to the ratio of this transverse momentum to the total momentum p . We get that

$$\alpha' = \frac{\Delta p_x}{p} = \frac{qwB}{p}. \quad (15.41)$$

We can compare this result with Eq. (15.39), which gives the same quantity computed quantum-mechanically. But the connection between classical mechanics and quantum mechanics is this: A particle of momentum p corresponds to a quan-

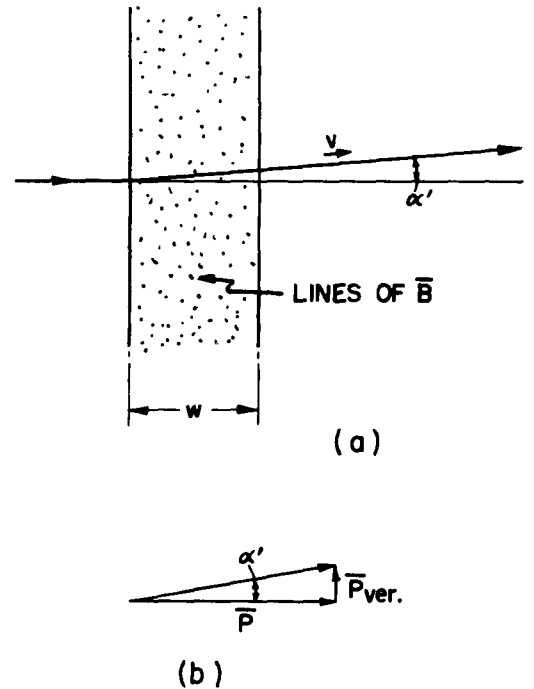


Fig. 15-9. Deflection of a particle due to passage through a strip of magnetic field.

tum amplitude varying with the wavelength $\lambda = \hbar/p$. With this equality, α and α' are identical; the classical and quantum calculations give the same result.

From the analysis we see how it is that the vector potential which appears in quantum mechanics in an explicit form produces a classical force which depends only on its derivatives. In quantum mechanics what matters is the interference between nearby paths; it always turns out that the effects depend only on how much the field A changes from point to point, and therefore only on the derivatives of A and not on the value itself. Nevertheless, the vector potential A (together with the scalar potential ϕ that goes with it) appears to give the most direct description of the physics. This becomes more and more apparent the more deeply we go into the quantum theory. In the general theory of quantum electrodynamics, one takes the vector and scalar potentials as the fundamental quantities in a set of equations that replace the Maxwell equations: E and B are slowly disappearing from the modern expression of physical laws; they are being replaced by A and ϕ .

15-6 What is true for statics is false for dynamics

We are now at the end of our exploration of the subject of static fields. Already in this chapter we have come perilously close to having to worry about what happens when fields change with time. We were barely able to avoid it in our treatment of magnetic energy by taking refuge in a relativistic argument. Even so, our treatment of the energy problem was somewhat artificial and perhaps even mysterious, because we ignored the fact that moving coils must, in fact, produce changing fields. It is now time to take up the treatment of time-varying fields—the subject of electrodynamics. We will do so in the next chapter. First, however, we would like to emphasize a few points.

Although we began this course with a presentation of the complete and correct equations of electromagnetism, we immediately began to study some incomplete pieces—because that was easier. There is a great advantage in starting with the simpler theory of static fields, and proceeding only later to the more complicated theory which includes dynamic fields. There is less new material to learn all at once, and there is time for you to develop your intellectual muscles in preparation for the bigger task.

But there is the danger in this process that before we get to see the complete story, the incomplete truths learned on the way may become ingrained and taken as the whole truth—that what is true and what is only sometimes true will become confused. So we give in Table 15-1 a summary of the important formulas we have covered, separating those which are true in general from those which are true for statics, but false for dynamics. This summary also shows, in part, where we are going, since as we treat dynamics we will be developing in detail what we must just state here without proof.

It may be useful to make a few remarks about the table. First, you should notice that the equations we started with are the *true* equations—we have not misled you there. The electromagnetic force (often called the *Lorentz force*) $\mathbf{F} = q(\mathbf{E} + \mathbf{v} \times \mathbf{B})$ is *true*. It is only Coulomb's law that is false, to be used only for statics. The four Maxwell equations for \mathbf{E} and \mathbf{B} are also true. The equations we took for statics are false, of course, because we left off all terms with time derivatives.

Gauss' law, $\nabla \cdot \mathbf{E} = \rho/\epsilon_0$, remains, but the curl of \mathbf{E} is *not* zero in general. So \mathbf{E} cannot always be equated to the gradient of a scalar—the electrostatic potential. We will see that a scalar potential still remains, but it is a time-varying quantity that must be used together with vector potentials for a complete description of the electric field. The equations governing this new scalar potential are, necessarily, also new.

We must also give up the idea that \mathbf{E} is zero in conductors. When the fields are changing, the charges in conductors do not, in general, have time to rearrange themselves to make the field zero. They are set in motion, but never reach equilibrium. The only general statement is: electric fields in conductors produce cur-

Table 15-1

FALSE IN GENERAL (true only for statics)	TRUE ALWAYS
$F = \frac{1}{4\pi\epsilon_0} \frac{q_1 q_2}{r^2} \quad (\text{Coulomb's law})$	$F = q(E + v \times B) \quad (\text{Lorentz force})$ $\rightarrow \nabla \cdot E = \frac{\rho}{\epsilon_0} \quad (\text{Gauss' law})$
$\nabla \times E = 0$ $E = -\nabla\phi$ $E(1) = \frac{1}{4\pi\epsilon_0} \frac{\rho(2)e_{12}}{r_{12}^2} dV_2$ <p>For conductors, $E = 0$, $\phi = \text{constant}$, $Q = CV$</p>	$\rightarrow \nabla \times E = -\frac{\partial B}{\partial t} \quad (\text{Faraday's law})$ $E = -\nabla\phi - \frac{\partial A}{\partial t}$ <p>In a conductor, E makes currents.</p>
$c^2 \nabla \times B = \frac{j}{\epsilon_0} \quad (\text{Ampere's law})$ $B(1) = \frac{1}{4\pi\epsilon_0 c^2} \int \frac{j(2) \times e_{12}}{r_{12}^2} dV_2$	$\rightarrow \nabla \cdot B = 0 \quad (\text{No magnetic charges})$ $B = \nabla \times A$ $\rightarrow c^2 \nabla \times B = \frac{j}{\epsilon_0} + \frac{\partial E}{\partial t}$
$\nabla^2 \phi = -\frac{\rho}{\epsilon_0} \quad (\text{Poisson's equation})$ $\left\{ \begin{array}{l} \nabla^2 A = -\frac{j}{\epsilon_0 c^2} \\ \text{with} \\ \nabla \cdot A = 0 \end{array} \right.$	$\left\{ \begin{array}{l} \nabla^2 \phi - \frac{1}{c^2} \frac{\partial^2 \phi}{\partial t^2} = -\frac{\rho}{\epsilon_0} \\ \text{and} \\ \nabla^2 A - \frac{1}{c^2} \frac{\partial^2 A}{\partial t^2} = -\frac{j}{\epsilon_0 c^2} \\ \text{with} \\ c^2 \nabla \cdot A + \frac{\partial \phi}{\partial t} = 0 \end{array} \right.$
$\phi(1) = \frac{1}{4\pi\epsilon_0} \int \frac{\rho(2)}{r_{12}} dV_2$ $A(1) = \frac{1}{4\pi\epsilon_0 c^2} \int \frac{j(2)}{r_{12}} dV_2$	$\left\{ \begin{array}{l} \phi(1, t) = \frac{1}{4\pi\epsilon_0} \int \frac{\rho(2, t')}{r_{12}} dV_2 \\ \text{and} \\ A(1, t) = \frac{1}{4\pi\epsilon_0 c^2} \int \frac{j(2, t')}{r_{12}} dV_2 \\ \text{with} \\ t' = t - \frac{r_{12}}{c} \end{array} \right.$
$U = \frac{1}{2} \int \rho \phi dV + \frac{1}{2} \int j \cdot A dV$	$U = \int \left(\frac{\epsilon_0}{2} E \cdot E + \frac{\epsilon_0 c^2}{2} B \cdot B \right) dV$

The equations marked by an arrow (\rightarrow) are Maxwell's equations.

rents. So in varying fields a conductor is *not* an equipotential. It also follows that the idea of a capacitance is no longer precise.

Since there are no magnetic charges, the divergence of \mathbf{B} is *always* zero. So \mathbf{B} can always be equated to $\nabla \times \mathbf{A}$. (Everything doesn't change!) But the generation of \mathbf{B} is not only from currents; $\nabla \times \mathbf{B}$ is proportional to the current density *plus* a new term $\partial \mathbf{E}/\partial t$. This means that \mathbf{A} is related to currents by a new equation. It is also related to ϕ . If we make use of our freedom to choose $\nabla \cdot \mathbf{A}$ for our own convenience, the equations for \mathbf{A} or ϕ can be arranged to take on a simple and elegant form. We therefore make the condition that $c^2 \nabla \cdot \mathbf{A} = -\partial \phi/\partial t$, and the differential equations for \mathbf{A} or ϕ appear as shown in the table.

The potentials \mathbf{A} and ϕ can still be found by integrals over the currents and charges, but not the *same* integrals as for statics. Most wonderfully, though, the true integrals are like the static ones, with only a small and physically appealing modification. When we do the integrals to find the potentials at some point, say point (1) in Fig. 15-10, we must use the values of \mathbf{j} and ρ at the point (2) *at an earlier time* $t' = t - r_{12}/c$. As you would expect, the influences propagate from point (2) to point (1) at the speed c . With this small change, one can solve for the fields of varying currents and charges, because once we have \mathbf{A} and ϕ , we get \mathbf{B} from $\nabla \times \mathbf{A}$, as before, and \mathbf{E} from $-\nabla \phi - \partial \mathbf{A}/\partial t$.

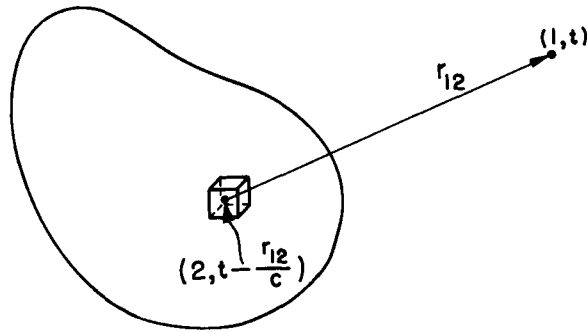


Fig. 15-10. The potentials at point (1) and at the time t are given by summing the contributions from each element of the source at the roving point (2), using the currents and charges which were present at the earlier time $t - r_{12}/c$.

Finally, you will notice that some results—for example, that the energy density in an electric field is $\epsilon_0 E^2/2$ —are true for electrodynamics as well as for statics. You should not be misled into thinking that this is at all “natural.” The validity of any formula derived in the static case must be demonstrated over again for the dynamic case. A contrary example is the expression for the electrostatic energy in terms of a volume integral of $\rho\phi$. This result is true *only* for statics.

We will consider all these matters in more detail in due time, but it will perhaps be useful to keep in mind this summary, so you will know what you can forget, and what you should remember as always true.

Induced Currents

16-1 Motors and generators

The discovery in 1820 that there was a close connection between electricity and magnetism was very exciting—until then, the two subjects had been considered as quite independent. The first discovery was that currents in wires make magnetic fields; then, in the same year, it was found that wires carrying current in a magnetic field have forces on them.

One of the excitements whenever there is a mechanical force is the possibility of using it in an engine to do work. Almost immediately after their discovery, people started to design electric motors using the forces on current-carrying wires. The principle of the electromagnetic motor is shown in bare outline in Fig. 16-1. A permanent magnet—usually with some pieces of soft iron—is used to produce a magnetic field in two slots. Across each slot there is a north and south pole, as shown. A rectangular coil of copper is placed with one side in each slot. When a current passes through the coil, it flows in opposite directions in the two slots, so the forces are also opposite, producing a torque on the coil about the axis shown. If the coil is mounted on a shaft so that it can turn, it can be coupled to pulleys or gears and can do work.

The same idea can be used for making a sensitive instrument for electrical measurements. Thus the moment the force law was discovered the precision of electrical measurements was greatly increased. First, the torque of such a motor can be made much greater for a given current by making the current go around many turns instead of just one. Then the coil can be mounted so that it turns with very little torque—either by supporting its shaft on very delicate jewel bearings or by hanging the coil on a very fine wire or a quartz fiber. Then an exceedingly small current will make the coil turn, and for small angles the amount of rotation will be proportional to the current. The rotation can be measured by gluing a pointer to the coil or, for the most delicate instruments, by attaching a small mirror to the coil and looking at the shift of the image of a scale. Such instruments are called galvanometers. Voltmeters and ammeters work on the same principle.

The same ideas can be applied on a large scale to make large motors for providing mechanical power. The coil can be made to go around and around by arranging that the connections to the coil are reversed each half-turn by contacts mounted on the shaft. Then the torque is always in the same direction. Small dc motors are made just this way. Larger motors, dc or ac, are often made by replacing the permanent magnet by an electromagnet, energized from the electrical power source.

With the realization that electric currents make magnetic fields, people immediately suggested that, somehow or other, magnets might also make electric fields. Various experiments were tried. For example, two wires were placed parallel to each other and a current was passed through one of them in the hope of finding a current in the other. The thought was that the magnetic field might in some way drag the electrons along in the second wire, giving some such law as “likes prefer to move alike.” With the largest available current and the most sensitive galvanometer to detect any current, the result was negative. Large magnets next to wires also produced no observed effects. Finally, Faraday discovered in 1840 the essential feature that had been missed—that electric effects exist only when there is something *changing*. If one of a pair of wires has a *changing* current, a current is induced in the other, or if a magnet is *moved* near an electric circuit, there is a current. We say that currents are *induced*. This was the induction effect discovered

16-1 Motors and generators

16-2 Transformers and inductances

16-3 Forces on induced currents

16-4 Electrical technology

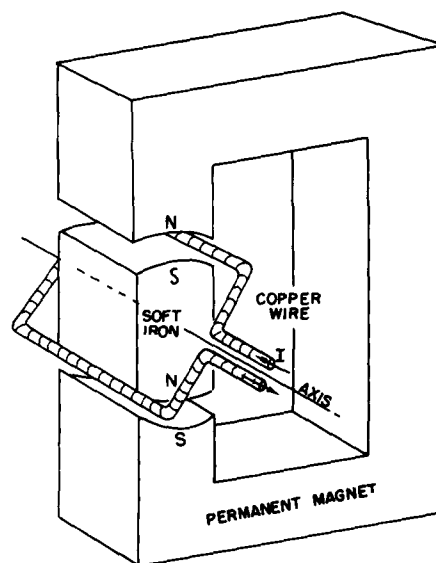


Fig. 16-1. Schematic outline of a simple electromagnetic motor.

by Faraday. It transformed the rather dull subject of static fields into a very exciting dynamic subject with an enormous range of wonderful phenomena. This chapter is devoted to a qualitative description of some of them. As we will see, one can quickly get into fairly complicated situations that are hard to analyze quantitatively in all their details. But never mind, our main purpose in this chapter is first to acquaint you with the phenomena involved. We will take up the detailed analysis later.

We can easily understand one feature of magnetic induction from what we already know, although it was not known in Faraday's time. It comes from the $\mathbf{v} \times \mathbf{B}$ force on a moving charge that is proportional to its velocity in a magnetic field. Suppose that we have a wire which passes near a magnet, as shown in Fig. 16-2, and that we connect the ends of the wire to a galvanometer. If we move the wire across the end of the magnet the galvanometer pointer moves.

The magnet produces some vertical magnetic field, and when we push the wire across the field, the electrons in the wire feel a *sideways* force—at right angles to the field and to the motion. The force pushes the electrons along the wire. But why does this move the galvanometer, which is so far from the force? Because when the electrons which feel the magnetic force try to move, they push—by electric repulsion—the electrons a little farther down the wire; they, in turn, repel the electrons a little farther on, and so on for a long distance. An amazing thing.

It was so amazing to Gauss and Weber—who first built a galvanometer—that they tried to see how far the forces in the wire would go. They strung a wire all the way across their city. Mr. Gauss, at one end, connected the wires to a battery (batteries were known before generators) and Mr. Weber watched the galvanometer move. They had a way of signaling long distances—it was the beginning of the telegraph! Of course, this has nothing directly to do with induction—it has to do with the way wires carry currents, whether the currents are pushed by induction or not.

Now suppose in the setup of Fig. 16-2 we leave the wire alone and move the magnet. We still see an effect on the galvanometer. As Faraday discovered, moving the magnet under the wire—one way—has the same effect as moving the wire over the magnet—the other way. But when the magnet is moved, we no longer have any $\mathbf{v} \times \mathbf{B}$ force on the electrons in the wire. This is the new effect that Faraday found. Today, we might hope to understand it from a relativity argument.

We already understand that the magnetic field of a magnet comes from its internal currents. So we expect to observe the same effect if instead of a magnet in Fig. 16-2 we use a coil of wire in which there is a current. If we move the wire past the coil there will be a current through the galvanometer, or also if we move the coil past the wire. But there is now a more exciting thing: If we change the magnetic field of the coil *not* by moving it, but by *changing its current*, there is again an effect in the galvanometer. For example, if we have a loop of wire near a coil, as shown in Fig. 16-3, and if we keep both of them stationary but switch off the current, there is a pulse of current through the galvanometer. When we switch the coil on again, the galvanometer kicks in the other direction.

Whenever the galvanometer in a situation such as the one shown in Fig. 16-2, or in Fig. 16-3, has a current, there is a net push on the electrons in the wire in one direction along the wire. There may be pushes in different directions at different places, but there is more push in one direction than another. What counts is the push integrated around the complete circuit. We call this net integrated push *the electromotive force* (abbreviated *emf*) in the circuit. More precisely, the emf is defined as the tangential force per unit charge in the wire integrated over length, once around the complete circuit. Faraday's complete discovery was that emf's can be generated in a wire in three different ways: by moving the wire, by moving a magnet near the wire, or by changing a current in a nearby wire.

Let's consider the simple machine of Fig. 16-1 again, only now, instead of putting a current through the wire to make it turn, let's turn the loop by an external force, for example by hand or by a waterwheel. When the coil rotates, its wires are moving in the magnetic field and we will find an emf in the circuit of the coil. The motor becomes a generator.

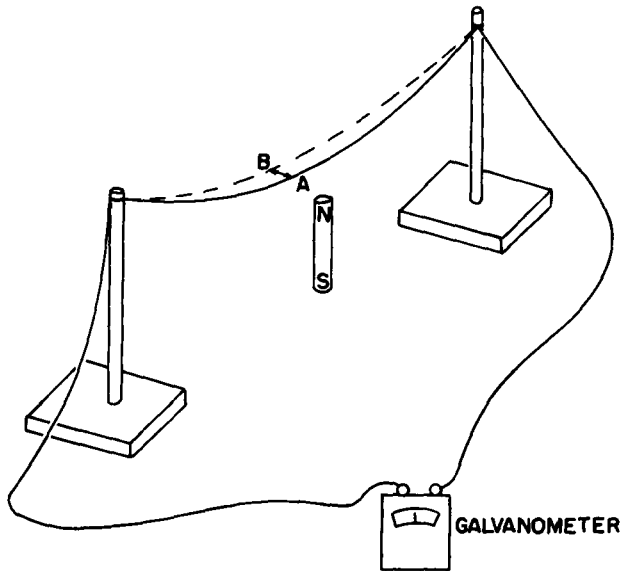


Fig. 16-2. Moving a wire through a magnetic field produces a current, as shown by the galvanometer.

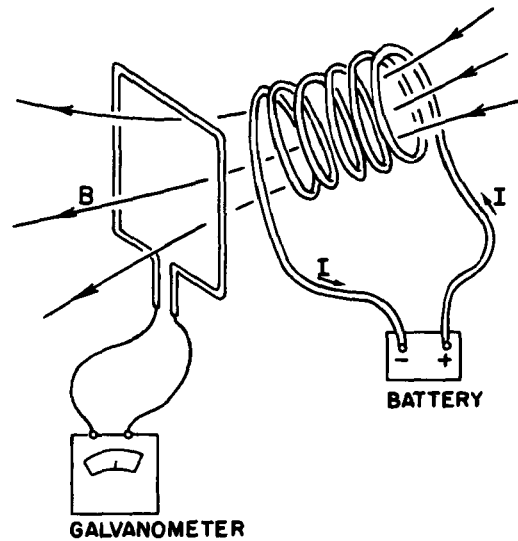


Fig. 16-3. A coil with current produces a current in a second coil if the first coil is moved or if its current is changed.

The coil of the generator has an induced emf from its motion. The amount of the emf is given by a simple rule discovered by Faraday. (We will just state the rule now and wait until later to examine it in detail.) The rule is that when the magnetic flux that passes through the loop (this flux is the normal component of B integrated over the area of the loop) is changing with time, the emf is equal to the rate of change of the flux. We will refer to this as "the flux rule." You see that when the coil of Fig. 16-1 is rotated, the flux through it changes. At the start some flux goes through one way; then when the coil has rotated 180° the same flux goes through the other way. If we continuously rotate the coil the flux is first positive, then negative, then positive, and so on. The rate of change of the flux must alternate also. So there is an alternating emf in the coil. If we connect the two ends of the coil to outside wires through some sliding contacts—called slip-rings—(just so the wires won't get twisted) we have an alternating-current generator.

Or we can also arrange, by means of some sliding contacts, that after every one-half rotation, the connection between the coil ends and the outside wires is reversed, so that when the emf reverses, so do the connections. Then the pulses of emf will always push currents in the same direction through the external circuit. We have what is called a direct-current generator.

The machine of Fig. 16-1 is either a motor or a generator. The reciprocity between motors and generators is nicely shown by using two identical dc "motors" of the permanent magnet kind, with their coils connected by two copper wires. When the shaft of one is turned mechanically, it becomes a generator and drives the other as a motor. If the shaft of the second is turned, it becomes the generator and drives the first as a motor. So here is an interesting example of a new kind of equivalence of nature: motor and generator are equivalent. The quantitative equivalence is, in fact, not completely accidental. It is related to the law of conservation of energy.

Another example of a device that can operate either to generate emf's or to respond to emf's is the receiver of a standard telephone—that is, an "earphone." The original telephone of Bell consisted of two such "earphones" connected by two long wires. The basic principle is shown in Fig. 16-4. A permanent magnet produces a magnetic field in two "yokes" of soft iron and in a thin diaphragm that is moved by sound pressure. When the diaphragm moves, it changes the amount of magnetic field in the yokes. Therefore a coil of wire wound around one of the yokes will have the flux through it changed when a sound wave hits the diaphragm.

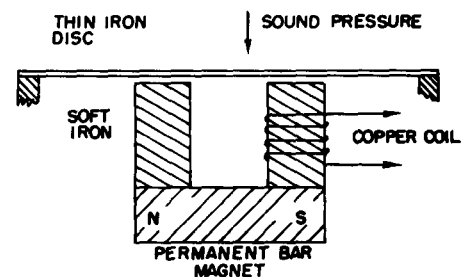


Fig. 16-4. A telephone transmitter or receiver.

So there is an emf in the coil. If the ends of the coil are connected to a circuit, a current which is an electrical representation of the sound is set up.

If the ends of the coil of Fig. 16-4 are connected by two wires to another identical gadget, varying currents will flow in the second coil. These currents will produce a varying magnetic field and will make a varying attraction on the iron diaphragm. The diaphragm will wiggle and make sound waves approximately similar to the ones that moved the original diaphragm. With a few bits of iron and copper the human voice is transmitted over wires!

(The modern home telephone uses a receiver like the one described but uses an improved invention to get a more powerful transmitter. It is the "carbon-button microphone," that uses sound pressure to vary the electric current from a battery.)

16-2 Transformers and inductances

One of the most interesting features of Faraday's discoveries is not that an emf exists in a moving coil—which we can understand in terms of the magnetic force $q\mathbf{v} \times \mathbf{B}$ —but that a changing current in one coil makes an emf in a second coil. And quite surprisingly the amount of emf induced in the second coil is given by the same "flux rule": that the emf is equal to the rate of change of the magnetic flux through the coil. Suppose that we take two coils, each wound around separate bundles of iron sheets (these help to make stronger magnetic fields), as shown in Fig. 16-5. Now we connect one of the coils—coil (a)—to an alternating-current generator. The continually changing current produces a continuously varying magnetic field. This varying field generates an alternating emf in the second coil—coil (b). This emf can, for example, produce enough power to light an electric bulb.

The emf alternates in coil (b) at a frequency which is, of course, the same as the frequency of the original generator. But the current in coil (b) can be larger or smaller than the current in coil (a). The current in coil (b) depends on the emf induced in it and on the resistance and inductance of the rest of its circuit. The emf can be less than that of the generator if, say, there is little flux change. Or the emf in coil (b) can be made much larger than that in the generator by winding coil (b) with many turns, since in a given magnetic field the flux through the coil is then greater. (Or if you prefer to look at it another way, the emf is the same in each turn, and since the total emf is the sum of the emf's of the separate turns, many turns in series produce a large emf.)

Such a combination of two coils—usually with an arrangement of iron sheets to guide the magnetic fields—is called a *transformer*. It can "transform" one emf (also called a "voltage") to another.

There are also induction effects in a single coil. For instance, in the setup in Fig. 16-5 there is a changing flux not only through coil (b), which lights the bulb, but also through coil (a). The varying current in coil (a) produces a varying magnetic field inside itself and the flux of this field is continually changing, so there is a *self-induced* emf in coil (a). There is an emf acting on any current when it is building up a magnetic field—or, in general, when its field is changing in any way. The effect is called *self-inductance*.

When we gave "the flux rule" that the emf is equal to the rate of change of the flux linkage, we didn't specify the direction of the emf. There is a simple rule, called Lenz's rule, for figuring out which way the emf goes: the emf *tries to oppose* any flux change. That is, the direction of an induced emf is always such that if a current were to flow in the direction of the emf, it would produce a flux of \mathbf{B} that opposes the change in \mathbf{B} that produces the emf. Lenz's rule can be used to find the direction of the emf in the generator of Fig. 16-1, or in the transformer winding of Fig. 16-3.

In particular, if there is a changing current in a single coil (or in any wire) there is a "back" emf in the circuit. This emf acts on the charges flowing in coil (a) of Fig. 16-5 to oppose the change in magnetic field, and so in the direction to oppose the change in current. It tries to keep the current constant; it is opposite to the current when the current is increasing, and it is in the direction of the current

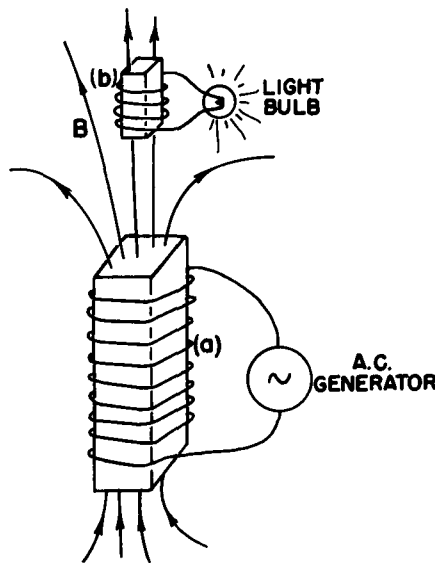


Fig. 16-5. Two coils, wrapped around bundles of iron sheets, allow a generator to light a bulb with no direct connection.

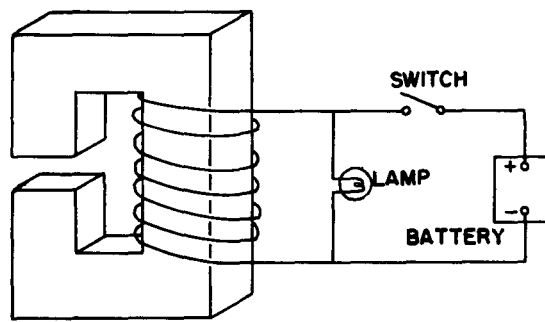


Fig. 16-6. Circuit connections for an electromagnet. The lamp allows the passage of current when the switch is opened, preventing the appearance of excessive emf's.

when it is decreasing. A current in a self-inductance has "inertia," because the inductive effects try to keep the flow constant, just as mechanical inertia tries to keep the velocity of an object constant.

Any large electromagnet will have a large self-inductance. Suppose that a battery is connected to the coil of a large electromagnet, as in Fig. 16-6, and that a strong magnetic field has been built up. (The current reaches a steady value determined by the battery voltage and the resistance of the wire in the coil.) But now suppose that we try to disconnect the battery by opening the switch. If we really opened the circuit, the current would go to zero rapidly, and in doing so it would generate an enormous emf. In most cases this emf would be large enough to develop an arc across the opening contacts of the switch. The high voltage that appears might also damage the insulation of the coil—or you, if you are the person who opens the switch! For these reasons, electromagnets are usually connected in a circuit like the one shown in Fig. 16-6. When the switch is opened, the current does not change rapidly but remains steady, flowing instead through the lamp, being driven by the emf from the self-inductance of the coil.

16-3 Forces on induced currents

You have probably seen the dramatic demonstration of Lenz's rule made with the gadget shown in Fig. 16-7. It is an electromagnet, just like coil (a) of Fig. 16-5. An aluminum ring is placed on the end of the magnet. When the coil is connected to an alternating-current generator by closing the switch, the ring flies into the air. The force comes, of course, from the induced currents in the ring. The fact that the ring flies away shows that the currents in it oppose the change of the field through it. When the magnet is making a north pole at its top, the induced current in the ring is making a downward-point north pole. The ring and the coil are repelled just like two magnets with like poles opposite. If a thin radial cut is made in the ring the force disappears, showing that it does indeed come from the currents in the ring.

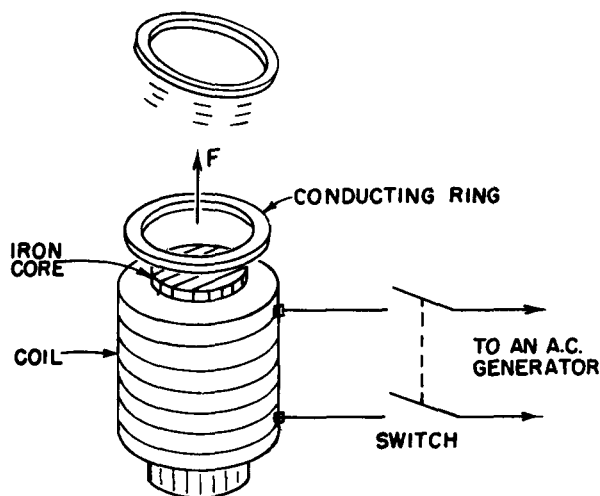


Fig. 16-7. A conducting ring is strongly repelled by an electromagnet with a varying current.

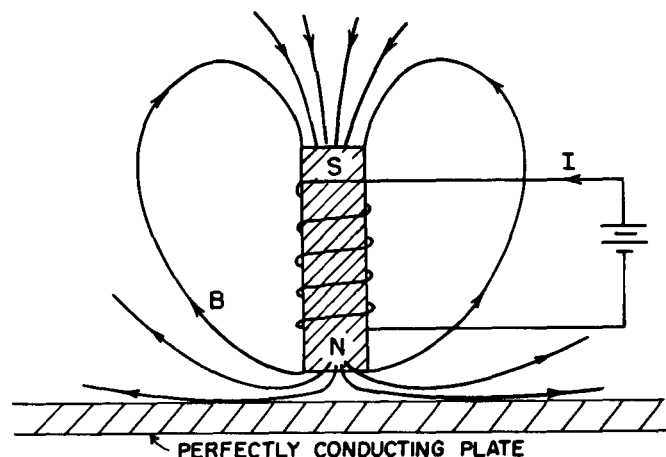


Fig. 16-8. An electromagnet near a perfectly conducting plate.

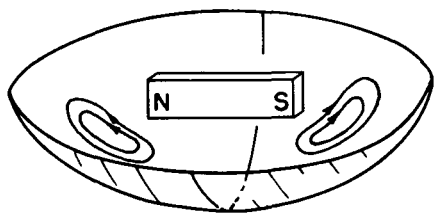


Fig. 16-9. A bar magnet is suspended above a superconducting bowl, by the repulsion of eddy currents.

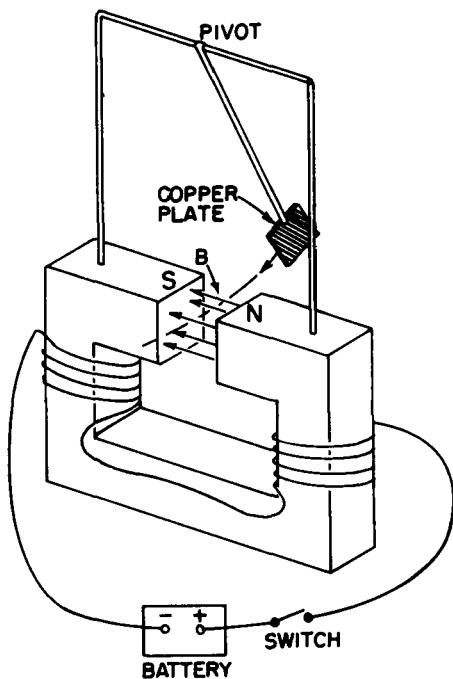


Fig. 16-10. The braking of the pendulum shows the forces due to eddy currents.

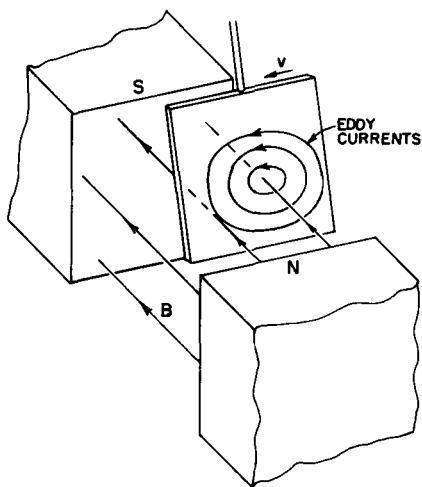


Fig. 16-11. The eddy currents in the copper pendulum.

If, instead of the ring, we place a disc of aluminum or copper across the end of the electromagnet of Fig. 16-7, it is also repelled; induced currents circulate in the material of the disc, and again produce a repulsion.

An interesting effect, similar in origin, occurs with a sheet of a perfect conductor. In a "perfect conductor" there is no resistance whatever to the current. So if currents are generated in it, they can keep going forever. In fact, the *slightest* emf would generate an arbitrarily large current—which really means that there can be no emf's at all. Any attempt to make a magnetic flux go through such a sheet generates currents that create opposite B fields—all with infinitesimal emf's, so with no flux entering.

If we have a sheet of a perfect conductor and put an electromagnet next to it, when we turn on the current in the magnet, currents called eddy currents appear in the sheet, so that no magnetic flux enters. The field lines would look as shown in Fig. 16-8. The same thing happens, of course, if we bring a bar magnet near a perfect conductor. Since the eddy currents are creating opposing fields, the magnets are repelled from the conductor. This makes it possible to suspend a bar magnet in air above a sheet of perfect conductor shaped like a dish, as shown in Fig. 16-9. The magnet is suspended by the repulsion of the induced eddy currents in the perfect conductor. There are no perfect conductors at ordinary temperatures, but some materials become perfect conductors at low enough temperatures. For instance, below 3.8°K tin conducts perfectly. It is called a superconductor.

If the conductor in Fig. 16-8 is not quite perfect there will be some resistance to flow of the eddy currents. The currents will tend to die out and the magnet will slowly settle down. The eddy currents in an imperfect conductor need an emf to keep them going, and to have an emf the flux must keep changing. The flux of the magnetic field gradually penetrates the conductor.

In a normal conductor, there are not only repulsive forces from eddy currents, but there can also be sidewise forces. For instance, if we move a magnet sideways along a conducting surface the eddy currents produce a force of drag, because the induced currents are opposing the changing of the location of flux. Such forces are proportional to the velocity and are like a kind of viscous force.

These effects show up nicely in the apparatus shown in Fig. 16-10. A square sheet of copper is suspended on the end of a rod to make a pendulum. The copper swings back and forth between the poles of an electromagnet. When the magnet is turned on, the pendulum motion is suddenly arrested. As the metal plate enters the gap of the magnet, there is a current induced in the plate which acts to oppose the change in flux through the plate. If the sheet were a perfect conductor, the currents would be so great that they would push the plate out again—it would bounce back. With a copper plate there is some resistance in the plate, so the currents at first bring the plate almost to a dead stop as it starts to enter the field. Then, as the currents die down, the plate slowly settles to rest in the magnetic field.

The nature of the eddy currents in the copper pendulum is shown in Fig. 16-11. The strength and geometry of the currents are quite sensitive to the shape of the plate. If, for instance, the copper plate is replaced by one which has several narrow slots cut in it, as shown in Fig. 16-12, the eddy-current effects are drastically reduced. The pendulum swings through the magnetic field with only a small retarding force. The reason is that the currents in each section of the copper have less flux to drive them, so the effects of the resistance of each loop are greater. The currents are smaller and the drag is less. The viscous character of the force is seen even more clearly if a sheet of copper is placed between the poles of the magnet of Fig. 16-10 and then released. It doesn't fall; it just sinks slowly downward. The eddy currents exert a strong resistance to the motion—just like the viscous drag in honey.

If, instead of dragging a conductor past a magnet, we try to rotate it in a magnetic field, there will be a resistive torque from the same effects. Alternatively, if we rotate a magnet—end over end—near a conducting plate or ring, the ring is dragged around; currents in the ring will create a torque that tends to rotate the ring with the magnet.

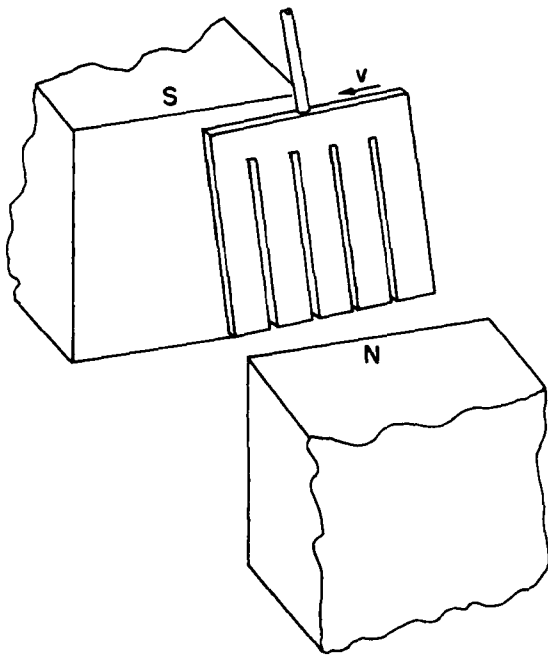


Fig. 16-12. Eddy-current effects are drastically reduced by cutting slots in the plate.

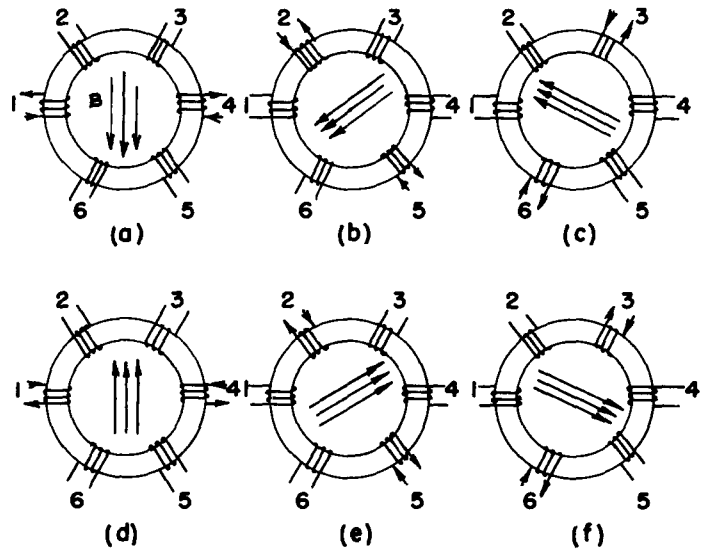


Fig. 16-13. Making a rotating magnetic field.

A field just like that of a rotating magnet can be made with an arrangement of coils such as is shown in Fig. 16-13. We take a torus of iron (that is, a ring of iron like a doughnut) and wind six coils on it. If we put a current, as shown in part (a), through windings (1) and (4), there will be a magnetic field in the direction shown in the figure. If we now switch the current to windings (2) and (5), the magnetic field will be in a new direction, as shown in part (b) of the figure. Continuing the process, we get the sequence of fields shown in the rest of the figure. If the process is done smoothly, we have a "rotating" magnetic field. We can easily get the required sequence of currents by connecting the coils to a three-phase power line, which provides just such a sequence of currents. "Three-phase power" is made in a generator using the principle of Fig. 16-1, except that there are *three* loops fastened together on the same shaft in a symmetrical way—that is, with an angle of 120° from one loop to the next. When the coils are rotated as a unit, the emf is a maximum in one, then in the next, and so on in a regular sequence. There are many practical advantages of three-phase power. One of them is the possibility of making a rotating magnetic field. The torque produced on a conductor by such a rotating field is easily shown by standing a metal ring on an insulating table just above the torus, as shown in Fig. 16-14. The rotating field causes the ring to spin about a vertical axis. The basic elements seen here are quite the same as those at play in a large commercial three-phase induction motor.

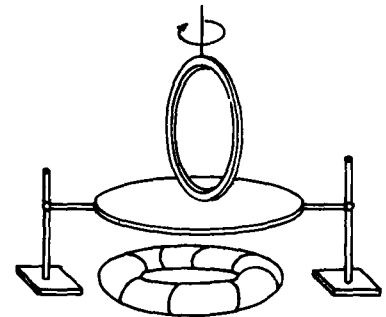


Fig. 16-14. The rotating field of Fig. 16-13 can be used to provide torque on a conducting ring.

Another form of induction motor is shown in Fig. 16-15. The arrangement shown is not suitable for a practical high-efficiency motor but will illustrate the principle. The electromagnet M , consisting of a bundle of laminated iron sheets wound with a solenoidal coil, is powered with alternating current from a generator. The magnet produces a varying flux of B through the aluminum disc. If we have just these two components, as shown in part (a) of the figure, we do not yet have a motor. There are eddy currents in the disc, but they are symmetric and there is no torque. (There will be some heating of the disc due to the induced currents.) If we now cover only one-half of the magnet pole with an aluminum plate, as shown in part (b) of the figure, the disc begins to rotate, and we have a motor. The operation depends on *two* eddy-current effects. First, the eddy currents in the aluminum plate oppose the change of flux through it, so the magnetic field above the plate always lags the field above that half of the pole which is not covered. This so-called "shaded-pole" effect produces a field which in the "shaded" region varies

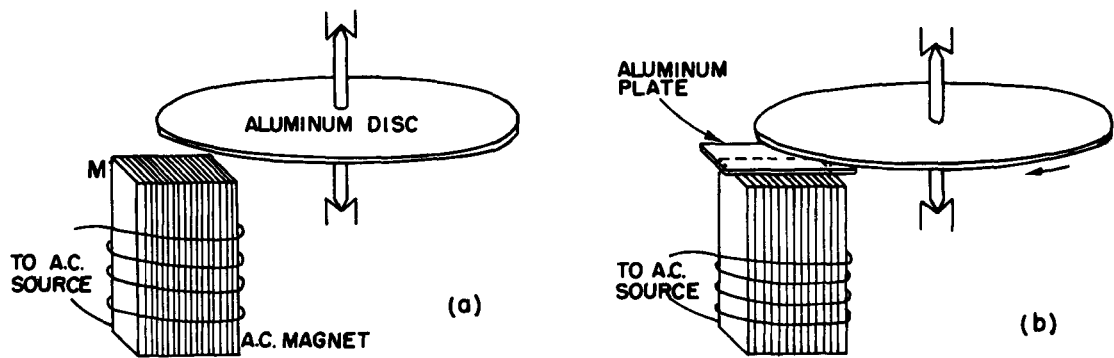


Fig. 16-15. A simple example of a shaded-pole induction motor.

much like that in the “unshaded” region except that it is delayed a constant amount in time. The whole effect is as if there were a magnet only half as wide which is continually being moved from the unshaded region toward the shaded one. Then the varying fields interact with the eddy currents in the disc to produce the torque on it.

16-4 Electrical technology

When Faraday first made public his remarkable discovery that a changing magnetic flux produces an emf, he was asked (as anyone is asked when he discovers a new fact of nature), “What is the use of it?” All he had found was the oddity that a tiny current was produced when he moved a wire near a magnet. Of what possible “use” could that be? His answer was: “What is the use of a new-born baby?”

Yet think of the tremendous practical applications his discovery has led to. What we have been describing are not just toys but examples chosen in most cases to represent the principle of some practical machine. For instance, the rotating ring in the turning field is an induction motor. There are, of course, some differences between it and a practical induction motor. The ring has a very small torque; it can be stopped with your hand. For a good motor, things have to be put together more intimately: there shouldn’t be so much “wasted” magnetic field out in the air. First, the field is concentrated by using iron. We have not discussed how iron does that, but iron can make the magnetic field tens of thousands of times stronger than copper coils alone could do. Second, the gaps between the pieces of iron are made small; to do that, some iron is even built into the rotating ring. Everything is arranged so as to get the greatest forces and the greatest efficiency—that is, conversion of electrical power to mechanical power—until the “ring” can no longer be held still by your hand.

This problem of closing the gaps and making the thing work in the most practical way is *engineering*. It requires serious study of design problems, although there are no new basic principles from which the forces are obtained. But there is a long way to go from the basic principles to a practical and economic design. Yet it is just such careful engineering design that has made possible such a tremendous thing as Boulder Dam and all that goes with it.

What is Boulder Dam? A huge river is stopped by a concrete wall. But what a wall it is! Shaped with a perfect curve that is very carefully worked out so that the least possible amount of concrete will hold back a whole river. It thickens at the bottom in that wonderful shape that the artists like but that the engineers can appreciate because they know that such thickening is related to the increase of pressure with the depth of the water. But we are getting away from electricity.

Then the water of the river is diverted into a huge pipe. That’s a nice engineering accomplishment in itself. The pipe feeds the water into a “waterwheel”—a huge turbine—and makes wheels turn. (Another engineering feat.) But why turn wheels? They are coupled to an exquisitely intricate mess of copper and iron, all

twisted and interwoven. With two parts—one that turns and one that doesn't. All a complex intermixture of a few materials, mostly iron and copper but also some paper and shellac for insulation. A revolving monster thing. A generator. Somewhere out of the mess of copper and iron come a few special pieces of copper. The dam, the turbine, the iron, the copper, all put there to make something special happen to a few bars of copper—an emf. Then the copper bars go a little way and circle for several times around another piece of iron in a transformer; then their job is done.

But around that same piece of iron curls another cable of copper which has no direct connection whatsoever to the bars from the generator; they have just been influenced because they passed near it—to get their emf. The transformer converts the power from the relatively low voltages required for the efficient design of the generator to the very high voltages that are best for efficient transmission of electrical energy over long cables.

And everything must be enormously efficient—there can be no waste, no loss. Why? The power for a metropolis is going through. If a small fraction were lost—one or two percent—think of the energy left behind! If one percent of the power were left in the transformer, that energy would need to be taken out somehow. If it appeared as heat, it would quickly melt the whole thing. There is, of course, some small inefficiency, but all that is required are a few pumps which circulate some oil through a radiator to keep the transformer from heating up.

Out of the Boulder Dam come a few dozen rods of copper—long, long, long rods of copper perhaps the thickness of your wrist that go for hundreds of miles in all directions. Small rods of copper carrying the power of a giant river. Then the rods are split to make more rods . . . then to more transformers . . . sometimes to great generators which recreate the current in another form . . . sometimes to engines turning for big industrial purposes . . . to more transformers . . . then more splitting and spreading . . . until finally the river is spread throughout the whole city—turning motors, making heat, making light, working gadgetry. The miracle of hot lights from cold water over 600 miles away—all done with specially arranged pieces of copper and iron. Large motors for rolling steel, or tiny motors for a dentist's drill. Thousands of little wheels, turning in response to the turning of the big wheel at Boulder Dam. Stop the big wheel, and all the wheels stop; the lights go out. They really are connected.

Yet there is more. The same phenomena that take the tremendous power of the river and spread it through the countryside, until a few drops of the river are running the dentist's drill, come again into the building of extremely fine instruments . . . for the detection of incredibly small amounts of current . . . for the transmission of voices, music, and pictures . . . for computers . . . for automatic machines of fantastic precision.

All this is possible because of carefully designed arrangements of copper and iron—efficiently created magnetic fields . . . blocks of rotating iron six feet in diameter whirling with clearances of 1/16 of an inch . . . careful proportions of copper for the optimum efficiency . . . strange shapes all serving a purpose, like the curve of the dam.

If some future archaeologist uncovers Boulder Dam, we may guess that he would admire the beauty of its curves. But also the explorers from some great future civilizations will look at the generators and transformers and say: "Notice that every iron piece has a beautifully efficient shape. Think of the thought that has gone into every piece of copper!"

This is the power of engineering and the careful design of our electrical technology. There has been created in the generator something which exists nowhere else in nature. It is true that there are forces of induction in other places. Certainly in some places around the sun and stars there are effects of electromagnetic induction. Perhaps also (though it's not certain) the magnetic field of the earth is maintained by an analog of an electric generator that operates on circulating currents in the interior of the earth. But nowhere have there been pieces put together with moving parts to generate electrical power as is done in the generator—with great efficiency and regularity.

You may think that designing electric generators is no longer an interesting subject, that it is a dead subject because they are all designed. Almost perfect generators or motors can be taken from a shelf. Even if this were true, we can admire the wonderful accomplishment of a problem solved to near perfection. But there remain as many unfinished problems. Even generators and transformers are returning as problems. It is likely that the whole field of low temperatures and superconductors will soon be applied to the problem of electric power distribution. With a radically new factor in the problem, new optimum designs will have to be created. Power networks of the future may have little resemblance to those of today.

You can see that there is an endless number of applications and problems that one could take up while studying the laws of induction. The study of the design of electrical machinery is a life work in itself. We cannot go very far in that direction, but we should be aware of the fact that when we have discovered the law of induction, we have suddenly connected our theory to an enormous practical development. We must, however, leave that subject to the engineers and applied scientists who are interested in working out the details of particular applications. Physics only supplies the base—the basic principles that apply, no matter what. (We have not yet completed the base, because we have yet to consider in detail the properties of iron and of copper. Physics has something to say about these as we will see a little later.)

Modern electrical technology began with Faraday's discoveries. The useless baby developed into a prodigy and changed the face of the earth in ways its proud father could never have imagined.

The Laws of Induction

17-1 The physics of induction

In the last chapter we described many phenomena which show that the effects of induction are quite complicated and interesting. Now we want to discuss the fundamental principles which govern these effects. We have already defined the emf in a conducting circuit as the total accumulated force on the charges throughout the length of the loop. More specifically, it is the tangential component of the force per unit charge, integrated along the wire once around the circuit. This quantity is equal, therefore, to the total work done on a single charge that travels once around the circuit.

We have also given the “flux rule,” which says that the emf is equal to the rate at which the magnetic flux through such a conducting circuit is changing. Let’s see if we can understand why that might be. First, we’ll consider a case in which the flux changes because a circuit is moved in a steady field.

In Fig. 17-1 we show a simple loop of wire whose dimensions can be changed. The loop has two parts, a fixed U-shaped part (a) and a movable crossbar (b) that can slide along the two legs of the U. There is always a complete circuit, but its area is variable. Suppose we now place the loop in a uniform magnetic field with the plane of the U perpendicular to the field. According to the rule, when the crossbar is moved there should be in the loop an emf that is proportional to the rate of change of the flux through the loop. This emf will cause a current in the loop. We will assume that there is enough resistance in the wire that the currents are small. Then we can neglect any magnetic field from this current.

The flux through the loop is wLB , so the “flux rule” would give for the emf—which we write as \mathcal{E} —

$$\mathcal{E} = wB \frac{dL}{dt} = wBv,$$

where v is the speed of translation of the crossbar.

Now we should be able to understand this result from the magnetic $\mathbf{v} \times \mathbf{B}$ forces on the charges in the moving crossbar. These charges will feel a force, tangential to the wire, equal to vB per unit charge. It is constant along the length w of the crossbar and zero elsewhere, so the integral is

$$\mathcal{E} = wvB,$$

which is the same result we got from the rate of change of the flux.

The argument just given can be extended to any case where there is a fixed magnetic field and the wires are moved. One can prove, in general, that for any circuit whose parts move in a fixed magnetic field the emf is the time derivative of the flux, regardless of the shape of the circuit.

On the other hand, what happens if the loop is stationary and the magnetic field is changed? We cannot deduce the answer to this question from the same argument. It was Faraday’s discovery—from experiment—that the “flux rule” is still correct no matter why the flux changes. The force on electric charges is given in complete generality by $\mathbf{F} = q(\mathbf{E} + \mathbf{v} \times \mathbf{B})$; there are no new special “forces due to changing magnetic fields.” Any forces on charges at rest in a stationary wire come from the \mathbf{E} term. Faraday’s observations led to the discovery that electric and magnetic fields are related by a new law: in a region where the magnetic field is changing with time, electric fields are generated. It is this electric

17-1 The physics of induction

17-2 Exceptions to the “flux rule”

17-3 Particle acceleration by an induced electric field; the betatron

17-4 A paradox

17-5 Alternating-current generator

17-6 Mutual inductance

17-7 Self-inductance

17-8 Inductance and magnetic energy

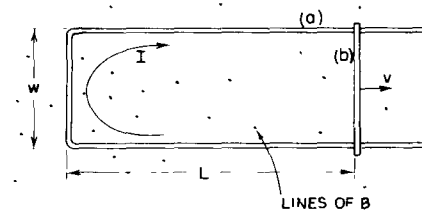


Fig. 17-1. An emf is induced in a loop if the flux is changed by varying the area of the circuit.

field which drives the electrons around the wire—and so is responsible for the emf in a stationary circuit when there is a changing magnetic flux.

The general law for the electric field associated with a changing magnetic field is

$$\nabla \times \mathbf{E} = -\frac{\partial \mathbf{B}}{\partial t}. \quad (17.1)$$

We will call this Faraday's law. It was discovered by Faraday but was first written in differential form by Maxwell, as one of his equations. Let's see how this equation gives the "flux rule" for circuits.

Using Stokes' theorem, this law can be written in integral form as

$$\oint_{\Gamma} \mathbf{E} \cdot d\mathbf{s} = \int_S (\nabla \times \mathbf{E}) \cdot \mathbf{n} \, da = - \int_S \frac{\partial \mathbf{B}}{\partial t} \cdot \mathbf{n} \, da, \quad (17.2)$$

where, as usual, Γ is any closed curve and S is any surface bounded by it. Here, remember, Γ is a *mathematical* curve fixed in space, and S is a fixed surface. Then the time derivative can be taken outside the integral and we have

$$\begin{aligned} \oint_{\Gamma} \mathbf{E} \cdot d\mathbf{s} &= -\frac{\partial}{\partial t} \int_S \mathbf{B} \cdot \mathbf{n} \, da \\ &= -\frac{\partial}{\partial t} (\text{flux through } S). \end{aligned} \quad (17.3)$$

Applying this relation to a curve Γ that follows a *fixed* circuit of conductor, we get the "flux rule" once again. The integral on the left is the emf, and that on the right is the negative rate of change of the flux linked by the circuit. So Eq. (17.1) applied to a fixed circuit is equivalent to the "flux rule."

So the "flux rule"—that the emf in a circuit is equal to the rate of change of the magnetic flux through the circuit—applies whether the flux changes because the field changes or because the circuit moves (or both). The two possibilities—"circuit moves" or "field changes"—are not distinguished in the statement of the rule. Yet in our explanation of the rule we have used two completely distinct laws for the two cases— $\mathbf{v} \times \mathbf{B}$ for "circuit moves" and $\nabla \times \mathbf{E} = -\partial \mathbf{B}/\partial t$ for "field changes."

We know of no other place in physics where such a simple and accurate general principle requires for its real understanding an analysis in terms of *two different phenomena*. Usually such a beautiful generalization is found to stem from a single deep underlying principle. Nevertheless, in this case there does not appear to be any such profound implication. We have to understand the "rule" as the combined effects of two quite separate phenomena.

We must look at the "flux rule" in the following way. In general, the force per unit charge is $\mathbf{F}/q = \mathbf{E} + \mathbf{v} \times \mathbf{B}$. In moving wires there is the force from the second term. Also, there is an \mathbf{E} -field if there is somewhere a changing magnetic field. They are independent effects, but the emf around the loop of wire is always equal to the rate of change of magnetic flux through it.

17-2 Exceptions to the "flux rule"

We will now give some examples, due in part to Faraday, which show the importance of keeping clearly in mind the distinction between the two effects responsible for induced emf's. Our examples involve situations to which the "flux rule" cannot be applied—either because there is no wire at all or because the *path* taken by induced currents moves about within an extended volume of a conductor.

We begin by making an important point: The part of the emf that comes from the \mathbf{E} -field does not depend on the existence of a physical wire (as does the $\mathbf{v} \times \mathbf{B}$ part). The \mathbf{E} -field can exist in free space, and its line integral around any imaginary line fixed in space is the rate of change of the flux of \mathbf{B} through that line. (Note that this is quite unlike the \mathbf{E} -field produced by static charges, for in that case the line integral of \mathbf{E} around a closed loop is always zero.)

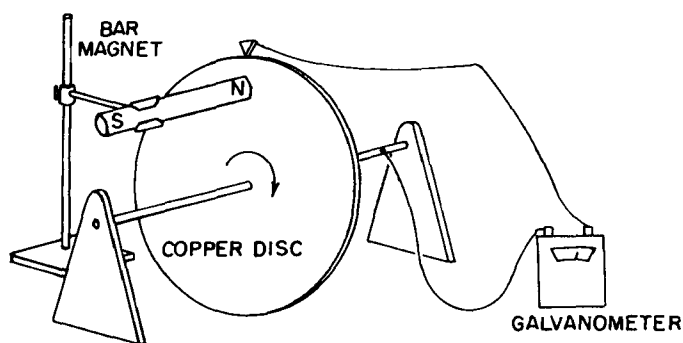


Fig. 17-2. When the disc rotates there is an emf from $\mathbf{v} \times \mathbf{B}$, but with no change in the linked flux.

Now we will describe a situation in which the flux through a circuit does not change, but there is nevertheless an emf. Figure 17-2 shows a conducting disc which can be rotated on a fixed axis in the presence of a magnetic field. One contact is made to the shaft and another rubs on the outer periphery of the disc. A circuit is completed through a galvanometer. As the disc rotates, the "circuit," in the sense of the place in space where the currents are, is always the same. But the part of the "circuit" in the disc is in material which is moving. Although the flux through the "circuit" is constant, there is still an emf, as can be observed by the deflection of the galvanometer. Clearly, here is a case where the $\mathbf{v} \times \mathbf{B}$ force in the moving disc gives rise to an emf which cannot be equated to a change of flux.

Now we consider, as an opposite example, a somewhat unusual situation in which the flux through a "circuit" (again in the sense of the place where the current is) changes but where there is *no* emf. Imagine two metal plates with slightly curved edges, as shown in Fig. 17-3, placed in a uniform magnetic field perpendicular to their surfaces. Each plate is connected to one of the terminals of a galvanometer, as shown. The plates make contact at one point P , so there is a complete circuit. If the plates are now rocked through a small angle, the point of contact will move to P' . If we imagine the "circuit" to be completed through the plates on the dotted line shown in the figure, the magnetic flux through this circuit changes by a large amount as the plates are rocked back and forth. Yet the rocking can be done with small motions, so that $\mathbf{v} \times \mathbf{B}$ is very small and there is practically no emf. The "flux rule" does not work in this case. It must be applied to circuits in which the *material* of the circuit remains the same. When the material of the circuit is changing, we must return to the basic laws. The *correct* physics is always given by the two basic laws

$$\mathbf{F} = q(\mathbf{E} + \mathbf{v} \times \mathbf{B}),$$

$$\nabla \times \mathbf{E} = -\frac{\partial \mathbf{B}}{\partial t}.$$

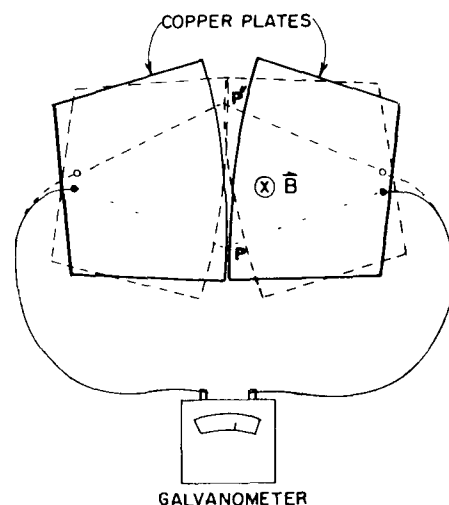


Fig. 17-3. When the plates are rocked in a uniform magnetic field, there can be a large change in the flux linkage without the generation of an emf.

17-3 Particle acceleration by an induced electric field; the betatron

We have said that the electromotive force generated by a changing magnetic field can exist even without conductors; that is, there can be magnetic induction without wires. We may still imagine an electromotive force around an arbitrary mathematical curve in space. It is defined as the tangential component of \mathbf{E} integrated around the curve. Faraday's law says that this line integral is equal to the rate of change of the magnetic flux through the closed curve, Eq. (17.3).

As an example of the effect of such an induced electric field, we want now to consider the motion of an electron in a changing magnetic field. We imagine a magnetic field which, everywhere on a plane, points in a vertical direction, as shown in Fig. 17-4. The magnetic field is produced by an electromagnet, but we will not worry about the details. For our example we will imagine that the magnetic field is symmetric about some axis, i.e., that the strength of the magnetic field will depend only on the distance from the axis. The magnetic field is also varying with time. We now imagine an electron that is moving in this field on a path that is a circle of constant radius with its center at the axis of the field. (We will see later

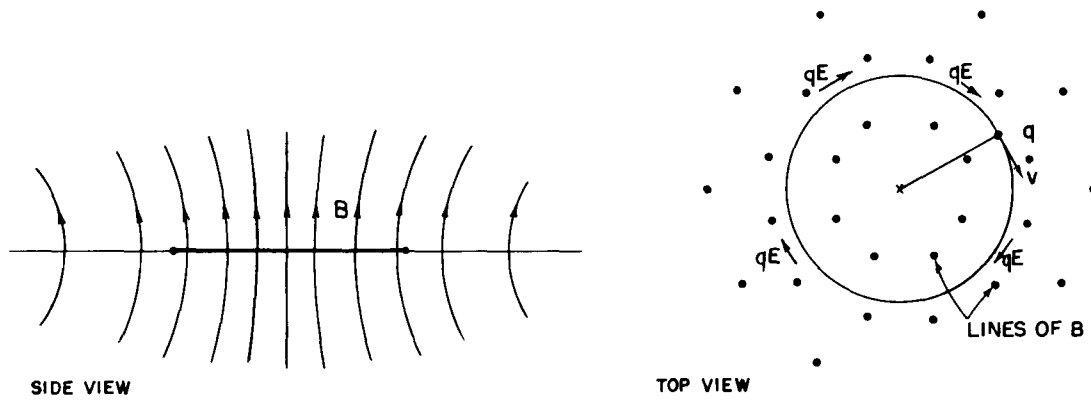


Fig. 17-4. An electron accelerating in an axially symmetric, time-varying magnetic field.

how this motion can be arranged.) Because of the changing magnetic field, there will be an electric field \mathbf{E} tangential to the electron's orbit which will drive it around the circle. Because of the symmetry, this electric field will have the same value everywhere on the circle. If the electron's orbit has the radius r , the line integral of \mathbf{E} around the orbit is equal to the rate of change of the magnetic flux through the circle. The line integral of \mathbf{E} is just its magnitude times the circumference of the circle, $2\pi r$. The magnetic flux must, in general, be obtained from an integral. For the moment, we let B_{av} represent the average magnetic field in the interior of the circle; then the flux is this average magnetic field times the area of the circle. We will have

$$2\pi r E = \frac{\partial}{\partial t} (B_{av} \cdot \pi r^2).$$

Since we are assuming r is constant, E is proportional to the time derivative of the average field:

$$E = \frac{r}{2} \frac{dB_{av}}{dt}. \quad (17.4)$$

The electron will feel the electric force $q\mathbf{E}$ and will be accelerated by it. Remembering that the relativistically correct equation of motion is that the rate of change of the momentum is proportional to the force, we have

$$qE = \frac{dp}{dt}. \quad (17.5)$$

For the circular orbit we have assumed, the electric force on the electron is always in the direction of its motion, so its total momentum will be increasing at the rate given by Eq. (17.5). Combining Eqs. (17.5) and (17.4), we may relate the rate of change of momentum to the change of the average magnetic field:

$$\frac{dp}{dt} = \frac{qr}{2} \frac{dB_{av}}{dt}. \quad (17.6)$$

Integrating with respect to t , we find for the electron's momentum

$$p = p_0 + \frac{qr}{2} \Delta B_{av}, \quad (17.7)$$

where p_0 is the momentum with which the electrons start out, and ΔB_{av} is the subsequent change in B_{av} . The operation of a *betatron*—a machine for accelerating electrons to high energies—is based on this idea.

To see how the betatron operates in detail, we must now examine how the electron can be constrained to move on a circle. We have discussed in Chapter 11 of Vol. I the principle involved. If we arrange that there is a magnetic field \mathbf{B} at the orbit of the electron, there will be a transverse force $q\mathbf{v} \times \mathbf{B}$ which, for a suit-

ably chosen \mathbf{B} , can cause the electron to keep moving on its assumed orbit. In the betatron this transverse force causes the electron to move in a circular orbit of constant radius. We can find out what the magnetic field at the orbit must be by using again the relativistic equation of motion, but this time, for the transverse component of the force. In the betatron (see Fig. 17-4), \mathbf{B} is at right angles to \mathbf{v} , so the transverse force is qvB . Thus the force is equal to the rate of change of the transverse component p_t of the momentum:

$$qvB = \frac{dp_t}{dt}. \quad (17.8)$$

When a particle is moving in a *circle*, the rate of change of its transverse momentum is equal to the magnitude of the total momentum times ω , the angular velocity of rotation (following the arguments of Chapter 11, Vol. I):

$$\frac{dp_t}{dt} = \omega p, \quad (17.9)$$

where, since the motion is circular,

$$\omega = \frac{v}{r}. \quad (17.10)$$

Setting the magnetic force equal to the transverse acceleration, we have

$$qvB_{\text{orbit}} = p \frac{v}{r}, \quad (17.11)$$

where B_{orbit} is the field at the radius r .

As the betatron operates, the momentum of the electron grows in proportion to B_{av} , according to Eq. (17.7), and if the electron is to continue to move in its proper circle, Eq. (17.11) must continue to hold as the momentum of the electron increases. The value of B_{orbit} must increase in proportion to the momentum p . Comparing Eq. (17.11) with Eq. (17.7), which determines p , we see that the following relation must hold between B_{av} , the average magnetic field *inside* the orbit at the radius r , and the magnetic field B_{orbit} at the orbit:

$$\Delta B_{\text{av}} = 2 \Delta B_{\text{orbit}}. \quad (17.12)$$

The correct operation of a betatron requires that the average magnetic field inside the orbit increase at twice the rate of the magnetic field at the orbit itself. In these circumstances, as the energy of the particle is increased by the induced electric field the magnetic field at the orbit increases at just the rate required to keep the particle moving in a circle.

The betatron is used to accelerate electrons to energies of tens of millions of volts, or even to hundreds of millions of volts. However, it becomes impractical for the acceleration of electrons to energies much higher than a few hundred million volts for several reasons. One of them is the practical difficulty of attaining the required high average value for the magnetic field inside the orbit. Another is that Eq. (17.6) is no longer correct at very high energies because it does not include the loss of energy from the particle due to its radiation of electromagnetic energy (the so-called synchrotron radiation discussed in Chapter 36, Vol. I). For these reasons, the acceleration of electrons to the highest energies—to many billions of electron volts—is accomplished by means of a different kind of machine, called a *synchrotron*.

17-4 A paradox

We would now like to describe for you an apparent paradox. A paradox is a situation which gives one answer when analyzed one way, and a different answer when analyzed another way, so that we are left in somewhat of a quandary as to actually what should happen. Of course, in physics there are never any real paradoxes because there is only one correct answer; at least we believe that nature will

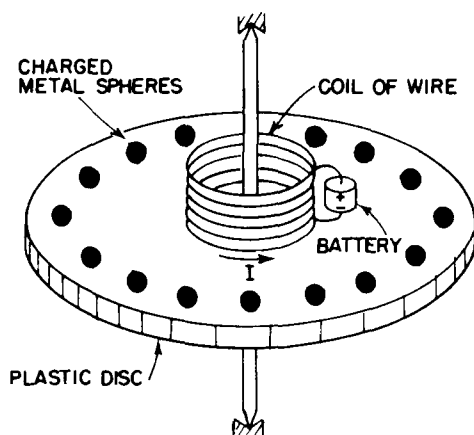


Fig. 17-5. Will the disc rotate if the current I is stopped?

act in only one way (and that is the *right way*, naturally). So in physics a paradox is only a confusion in our own understanding. Here is our paradox.

Imagine that we construct a device like that shown in Fig. 17-5. There is a thin, circular plastic disc supported on a concentric shaft with excellent bearings, so that it is quite free to rotate. On the disc is a coil of wire in the form of a short solenoid concentric with the axis of rotation. This solenoid carries a steady current I provided by a small battery, also mounted on the disc. Near the edge of the disc and spaced uniformly around its circumference are a number of small metal spheres insulated from each other and from the solenoid by the plastic material of the disc. Each of these small conducting spheres is charged with the same electrostatic charge Q . Everything is quite stationary, and the disc is at rest. Suppose now that by some accident—or by prearrangement—the current in the solenoid is interrupted, without, however, any intervention from the outside. So long as the current continued, there was a magnetic flux through the solenoid more or less parallel to the axis of the disc. When the current is interrupted, this flux must go to zero. There will, therefore, be an electric field induced which will circulate around in circles centered at the axis. The charged spheres on the perimeter of the disc will all experience an electric field tangential to the perimeter of the disc. This electric force is in the same sense for all the charges and so will result in a net torque on the disc. From these arguments we would expect that as the current in the solenoid disappears, the disc would begin to rotate. If we knew the moment of inertia of the disc, the current in the solenoid, and the charges on the small spheres, we could compute the resulting angular velocity.

But we could also make a different argument. Using the principle of the conservation of angular momentum, we could say that the angular momentum of the disc with all its equipment is initially zero, and so the angular momentum of the assembly should remain zero. There should be no rotation when the current is stopped. Which argument is correct? Will the disc rotate or will it not? We will leave this question for you to think about.

We should warn you that the correct answer does not depend on any non-essential feature, such as the asymmetric position of a battery, for example. In fact, you can imagine an ideal situation such as the following. The solenoid is made of superconducting wire through which there is a current. After the disc has been carefully placed at rest, the temperature of the solenoid is allowed to rise slowly. When the temperature of the wire reaches the transition temperature between superconductivity and normal conductivity, the current in the solenoid will be brought to zero by the resistance of the wire. The flux will, as before, fall to zero, and there will be an electric field around the axis. We should also warn you that the solution is not easy, nor is it a trick. When you figure it out, you will have discovered an important principle of electromagnetism.

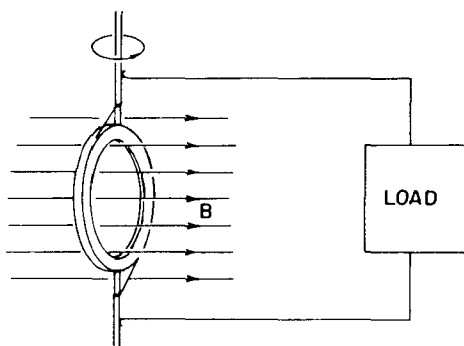


Fig. 17-6. A coil of wire rotating in a uniform magnetic field—the basic idea of the ac generator.

17-5 Alternating-current generator

In the remainder of this chapter we apply the principles of Section 17-1 to analyze a number of the phenomena discussed in Chapter 16. We first look in more detail at the alternating-current generator. Such a generator consists basically of a coil of wire rotating in a uniform magnetic field. The same result can also be achieved by a fixed coil in a magnetic field whose direction rotates in the manner described in the last chapter. We will consider only the former case. Suppose we have a circular coil of wire which can be turned on an axis along one of its diameters. Let this coil be located in a uniform magnetic field perpendicular to the axis of rotation, as in Fig. 17-6. We also imagine that the two ends of the coil are brought to external connections through some kind of sliding contacts.

Due to the rotation of the coil, the magnetic flux through it will be changing. The circuit of the coil will therefore have an emf in it. Let S be the area of the coil and θ the angle between the magnetic field and the normal to the plane of the coil.*

* Now that we are using the letter A for the vector potential, we prefer to let S stand for a Surface area.

The flux through the coil is then

$$BS \cos \theta. \quad (17.13)$$

If the coil is rotating at the uniform angular velocity ω , θ varies with time as $\theta = \omega t$. The emf \mathcal{E} in the coil is then

$$\mathcal{E} = -\frac{d}{dt}(\text{flux}) = -\frac{d}{dt}(BS \cos \omega t),$$

or

$$\mathcal{E} = BS\omega \sin \omega t. \quad (17.14)$$

If we bring the wires from the generator to a point some distance from the rotating coil, where the magnetic field is zero, or at least is not varying with time, the curl of \mathbf{E} in this region will be zero and we can define an electric potential. In fact, if there is no current being drawn from the generator, the potential difference V between the two wires will be equal to the emf in the rotating coil. That is,

$$V = BS\omega \sin \omega t = V_0 \sin \omega t.$$

The potential difference between the wires varies as $\sin \omega t$. Such a varying potential difference is called an alternating voltage.

Since there is an electric field between the wires, they must be electrically charged. It is clear that the emf of the generator has pushed some excess charges out to the wire until the electric field from them is strong enough to exactly counterbalance the induction force. Seen from outside the generator, the two wires appear as though they had been electrostatically charged to the potential difference V , and as though the charge was being changed with time to give an alternating potential difference. There is also another difference from an electrostatic situation. If we connect the generator to an external circuit that permits passage of a current, we find that the emf does not permit the wires to be discharged but continues to provide charge to the wires as current is drawn from them, attempting to keep the wires always at the same potential difference. If, in fact, the generator is connected in a circuit whose total resistance is R , the current through the circuit will be proportional to the emf of the generator and inversely proportional to R . Since the emf has a sinusoidal time variation, so also does the current. There is an alternating current

$$I = \frac{\mathcal{E}}{R} = \frac{V_0}{R} \sin \omega t.$$

The schematic diagram of such a circuit is shown in Fig. 17-7.

We can also see that the emf determines how much energy is supplied by the generator. Each charge in the wire is receiving energy at the rate $\mathbf{F} \cdot \mathbf{v}$, where \mathbf{F} is the force on the charge and \mathbf{v} is its velocity. Now let the number of moving charges per unit length of the wire be n ; then the power being delivered into any element ds of the wire is

$$\mathbf{F} \cdot \mathbf{v} n ds.$$

For a wire, \mathbf{v} is always along $d\mathbf{s}$, so we can rewrite the power as

$$n\mathbf{F} \cdot d\mathbf{s}.$$

The total power being delivered to the complete circuit is the integral of this expression around the complete loop:

$$\text{Power} = \oint n\mathbf{F} \cdot d\mathbf{s}. \quad (17.15)$$

Now remember that $qn\mathbf{v}$ is the current \mathbf{I} , and that the emf is defined as the integral of \mathbf{F}/q around the circuit. We get the result

$$\text{Power from a generator} = \mathcal{E}I. \quad (17.16)$$

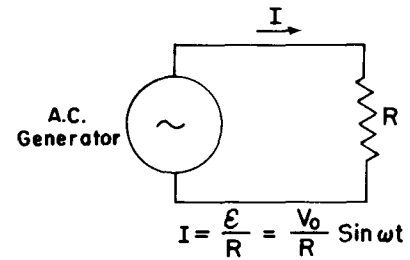


Fig. 17-7. A circuit with an ac generator and a resistance.

When there is a current in the coil of the generator, there will also be mechanical forces on it. In fact, we know that the torque on the coil is proportional to its magnetic moment, to the magnetic field strength B , and to the sine of the angle between. The magnetic moment is the current in the coil times its area. Therefore the torque is

$$\tau = ISB \sin \theta. \quad (17.17)$$

The rate at which mechanical work must be done to keep the coil rotating is the angular velocity ω times the torque:

$$\frac{dW}{dt} = \omega\tau = \omega ISB \sin \theta. \quad (17.18)$$

Comparing this equation with Eq. (17.14), we see that the rate of mechanical work required to rotate the coil against the magnetic forces is just equal to $\mathcal{E}I$, the rate at which electrical energy is delivered by the emf of the generator. All of the mechanical energy used up in the generator appears as electrical energy in the circuit.

As another example of the currents and forces due to an induced emf, let's analyze what happens in the setup described in Section 12, and shown in Fig. 17-1. There are two parallel wires and a sliding crossbar located in a uniform magnetic field perpendicular to the plane of the parallel wires. Now let's assume that the "bottom" of the U (the left side in the figure) is made of wires of high resistance, while the two side wires are made of a good conductor like copper—then we don't need to worry about the change of the circuit resistance as the crossbar is moved. As before, the emf in the circuit is

$$\mathcal{E} = vBw. \quad (17.19)$$

The current in the circuit is proportional to this emf and inversely proportional to the resistance of the circuit:

$$I = \frac{\mathcal{E}}{R} = \frac{vBw}{R}. \quad (17.20)$$

Because of this current there will be a magnetic force on the crossbar that is proportional to its length, to the current in it, and to the magnetic field, such that

$$F = BIw. \quad (17.21)$$

Taking I from Eq. (17.20), we have for the force

$$F = \frac{B^2 w^2}{R} v. \quad (17.22)$$

We see that the force is proportional to the velocity of the crossbar. The direction of the force, as you can easily see, is opposite to its velocity. Such a "velocity-proportional" force, which is like the force of viscosity, is found whenever induced currents are produced by moving conductors in a magnetic field. The examples of eddy currents we gave in the last chapter also produced forces on the conductors proportional to the velocity of the conductor, even though such situations, in general, give a complicated distribution of currents which is difficult to analyze.

It is often convenient in the design of mechanical systems to have damping forces which are proportional to the velocity. Eddy-current forces provide one of the most convenient ways of getting such a velocity-dependent force. An example of the application of such a force is found in the conventional domestic wattmeter. In the wattmeter there is a thin aluminum disc that rotates between the poles of a permanent magnet. This disc is driven by a small electric motor whose torque is proportional to the power being consumed in the electrical circuit of the house. Because of the eddy-current forces in the disc, there is a resistive force proportional to the velocity. In equilibrium, the velocity is therefore proportional to the rate of consumption of electrical energy. By means of a counter attached to the rotating disc, a record is kept of the number of revolutions it makes. This count is an indication of the total energy consumption, i.e., the number of watthours used.

We may also point out that Eq. (17.22) shows that the force from induced currents—that is, any eddy-current force—is inversely proportional to the resistance. The force will be larger, the better the conductivity of the material. The reason, of course, is that an emf produces more current if the resistance is low, and the stronger currents represent greater mechanical forces.

We can also see from our formulas how mechanical energy is converted into electrical energy. As before, the electrical energy supplied to the resistance of the circuit is the product $\mathcal{E}I$. The rate at which work is done in moving the conducting crossbar is the force on the bar times its velocity. Using Eq. (17.21) for the force, the rate of doing work is

$$\frac{dW}{dt} = \frac{v^2 B^2 w^2}{R}.$$

We see that this is indeed equal to the product $\mathcal{E}I$ we would get from Eqs. (17.19) and (17.20). Again the mechanical work appears as electrical energy.

17-6 Mutual inductance

We now want to consider a situation in which there are fixed coils of wire but changing magnetic fields. When we described the production of magnetic fields by currents, we considered only the case of steady currents. But so long as the currents are changed slowly, the magnetic field will at each instant be nearly the same as the magnetic field of a steady current. We will assume in the discussion of this section that the currents are always varying sufficiently slowly that this is true.

In Fig. 17-8 is shown an arrangement of two coils which demonstrates the basic effects responsible for the operation of a transformer. Coil 1 consists of a conducting wire wound in the form of a long solenoid. Around this coil—and insulated from it—is wound coil 2, consisting of a few turns of wire. If now a current is passed through coil 1, we know that a magnetic field will appear inside it. This magnetic field also passes through coil 2. As the current in coil 1 is varied, the magnetic flux will also vary, and there will be an induced emf in coil 2. We will now calculate this induced emf.

We have seen in Section 13-5 that the magnetic field inside a long solenoid is uniform and has the magnitude

$$B = \frac{1}{\epsilon_0 c^2} \frac{N_1 I_1}{l}, \quad (17.23)$$

where N_1 is the number of turns in coil 1, I_1 is the current through it, and l is its length. Let's say that the cross-sectional area of coil 1 is S ; then the flux of \mathbf{B} is its magnitude times S . If coil 2 has N_2 turns, this flux links the coil N_2 times. Therefore the emf in coil 2 is given by

$$\mathcal{E}_2 = -N_2 S \frac{dB}{dt}. \quad (17.24)$$

The only quantity in Eq. (17.23) which varies with time is I_1 . The emf is therefore given by

$$\mathcal{E}_2 = -\frac{N_1 N_2 S}{\epsilon_0 c^2 l} \frac{dI_1}{dt}. \quad (17.25)$$

We see that the emf in coil 2 is proportional to the rate of change of the current in coil 1. The constant of proportionality, which is basically a geometric factor of the two coils, is called the *mutual inductance*, and is usually designated \mathfrak{M}_{21} . Equation (17.25) is then written

$$\mathcal{E}_2 = \mathfrak{M}_{21} \frac{dI_1}{dt}. \quad (17.26)$$

Suppose now that we were to pass a current through coil 2 and ask about the emf in coil 1. We would compute the magnetic field, which is everywhere

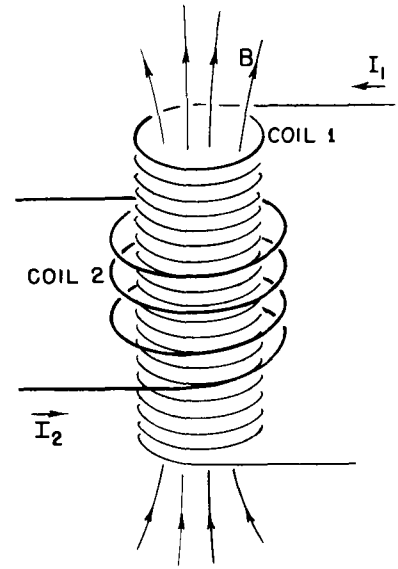


Fig. 17-8. A current in coil 1 produces a magnetic field through coil 2.

proportional to the current I_2 . The flux linkage through coil 1 would depend on the geometry, but would be proportional to the current I_2 . The emf in coil 1 would, therefore, again be proportional to dI_2/dt : We can write

$$\varepsilon_1 = \mathfrak{M}_{12} \frac{dI_2}{dt}. \quad (17.27)$$

The computation of \mathfrak{M}_{12} would be more difficult than the computation we have just done for \mathfrak{M}_{21} . We will not carry through that computation now, because we will show later in this chapter that \mathfrak{M}_{12} is necessarily equal to \mathfrak{M}_{21} .

Since for *any* coil its field is proportional to its current, the same kind of result would be obtained for any two coils of wire. The equations (17.26) and (17.27) would have the same form; only the constants \mathfrak{M}_{21} and \mathfrak{M}_{12} would be different. Their values would depend on the shapes of the coils and their relative positions.

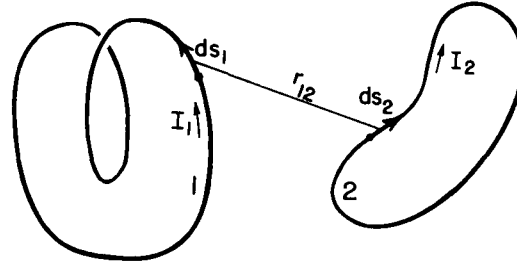


Fig. 17-9. Any two coils have a mutual inductance \mathfrak{M} proportional to the integral of $ds_1 \cdot ds_2/r_{12}$.

Suppose that we wish to find the mutual inductance between any two arbitrary coils—for example, those shown in Fig. 17-9. We know that the general expression for the emf in coil 1 can be written as

$$\varepsilon_1 = -\frac{d}{dt} \int_{(1)} \mathbf{B} \cdot \mathbf{n} \, da,$$

where \mathbf{B} is the magnetic field and the integral is to be taken over a surface bounded by circuit 1. We have seen in Section 14-1 that such a surface integral of \mathbf{B} can be related to a line integral of the vector potential. In particular,

$$\int_{(1)} \mathbf{B} \cdot \mathbf{n} \, da = \oint_{(1)} \mathbf{A} \cdot d\mathbf{s}_1,$$

where \mathbf{A} represents the vector potential and $d\mathbf{s}_1$ is an element of circuit 1. The line integral is to be taken around circuit 1. The emf in coil 1 can therefore be written as

$$\varepsilon_1 = -\frac{d}{dt} \oint_{(1)} \mathbf{A} \cdot d\mathbf{s}_1. \quad (17.28)$$

Now let's assume that the vector potential at circuit 1 comes from currents in circuit 2. Then it can be written as a line integral around circuit 2:

$$\mathbf{A} = \frac{1}{4\pi\epsilon_0 c^2} \oint_{(2)} \frac{I_2 d\mathbf{s}_2}{r_{12}}, \quad (17.29)$$

where I_2 is the current in circuit 2, and r_{12} is the distance from the element of the circuit $d\mathbf{s}_2$ to the point on circuit 1 at which we are evaluating the vector potential. (See Fig. 17-9.) Combining Eqs. (17.28) and (17.29), we can express the emf in circuit 1 as a double line integral:

$$\varepsilon_1 = -\frac{1}{4\pi\epsilon_0 c^2} \frac{d}{dt} \oint_{(1)} \oint_{(2)} \frac{I_2 d\mathbf{s}_2}{r_{12}} \cdot d\mathbf{s}_1.$$

In this equation the integrals are all taken with respect to stationary circuits. The only variable quantity is the current I_2 , which does not depend on the variables of

integration. We may therefore take it out of the integrals. The emf can then be written as

$$\mathcal{E}_1 = \mathfrak{M}_{12} \frac{dI_2}{dt},$$

where the coefficient \mathfrak{M}_{12} is

$$\mathfrak{M}_{12} = -\frac{1}{4\pi\epsilon_0 c^2} \oint_{(1)} \oint_{(2)} \frac{ds_2 \cdot ds_1}{r_{12}}. \quad (17.30)$$

We see from this integral that \mathfrak{M}_{12} depends only on the circuit geometry. It depends on a kind of average separation of the two circuits, with the average weighted most for parallel segments of the two coils. Our equation can be used for calculating the mutual inductance of any two circuits of arbitrary shape. Also, it shows that the integral for \mathfrak{M}_{12} is identical to the integral for \mathfrak{M}_{21} . We have therefore shown that the two coefficients are identical. For a system with only two coils, the coefficients \mathfrak{M}_{12} and \mathfrak{M}_{21} are often represented by the symbol \mathfrak{M} without subscripts, called simply the *mutual inductance*:

$$\mathfrak{M}_{12} = \mathfrak{M}_{21} = \mathfrak{M}.$$

17-7 Self-inductance

In discussing the induced electromotive forces in the two coils of Figs. 17-8 or 17-9, we have considered only the case in which there was a current in one coil or the other. If there are currents in the two coils simultaneously, the magnetic flux linking either coil will be the sum of the two fluxes which would exist separately, because the law of superposition applies for magnetic fields. The emf in either coil will therefore be proportional not only to the change of the current in the other coil, but also to the change in the current of the coil itself. Thus the total emf in coil 2 should be written*

$$\mathcal{E}_2 = \mathfrak{M}_{21} \frac{dI_1}{dt} + \mathfrak{M}_{22} \frac{dI_2}{dt}. \quad (17.31)$$

Similarly, the emf in coil 1 will depend not only on the changing current in coil 2, but also on the changing current in itself:

$$\mathcal{E}_1 = \mathfrak{M}_{12} \frac{dI_2}{dt} + \mathfrak{M}_{11} \frac{dI_1}{dt}. \quad (17.32)$$

The coefficients \mathfrak{M}_{22} and \mathfrak{M}_{11} are always negative numbers. It is usual to write

$$\mathfrak{M}_{11} = -\mathcal{L}_1, \quad \mathfrak{M}_{22} = -\mathcal{L}_2, \quad (17.33)$$

where \mathcal{L}_1 and \mathcal{L}_2 are called the *self-inductances* of the two coils.

The self-induced emf will, of course, exist even if we have only one coil. Any coil by itself will have a self-inductance \mathcal{L} . The emf will be proportional to the rate of change of the current in it. For a single coil, it is usual to adopt the convention that the emf and the current are considered positive if they are in the same direction. With this convention, we may write for the emf of a single coil

$$\mathcal{E} = -\mathcal{L} \frac{dI}{dt}. \quad (17.34)$$

The negative sign indicates that the emf opposes the change in current—it is often called a “back emf.”

Since any coil has a self-inductance which opposes the change in current, the current in the coil has a kind of inertia. In fact, if we wish to change the current in

* The sign of \mathfrak{M}_{12} and \mathfrak{M}_{21} in Eqs. (17.31) and (17.32) depends on the arbitrary choices for the sense of a positive current in the two coils.

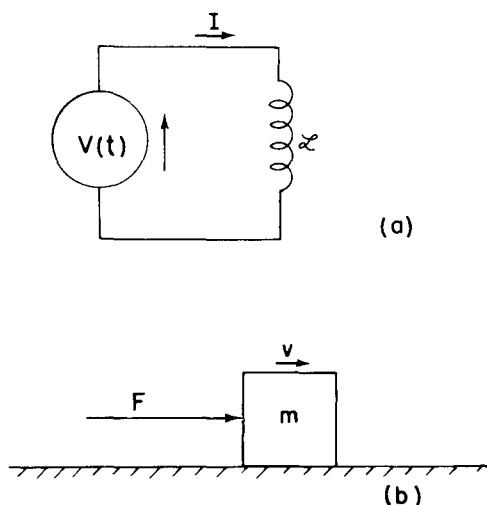


Fig. 17-10 (a) A circuit with a voltage source and an inductance. (b) An analogous mechanical system.

a coil we must overcome this inertia by connecting the coil to some external voltage source such as a battery or a generator, as shown in the schematic diagram of Fig. 17-10(a). In such a circuit, the current I depends on the voltage \mathcal{V} according to the relation

$$\mathcal{V} = \mathcal{L} \frac{dI}{dt}. \quad (17.35)$$

This equation has the same form as Newton's law of motion for a particle in one dimension. We can therefore study it by the principle that "the same equations have the same solutions." Thus, if we make the externally applied voltage \mathcal{V} correspond to an externally applied force F , and the current I in a coil correspond to the velocity v of a particle, the inductance \mathcal{L} of the coil corresponds to the mass m of the particle.* See Fig. 17-10(b). We can make the following table of corresponding quantities.

Particle	Coil
F (force)	\mathcal{V} (potential difference)
v (velocity)	I (current)
x (displacement)	q (charge)
$F = m \frac{dv}{dt}$	$\mathcal{V} = \mathcal{L} \frac{dI}{dt}$
mv (momentum)	$\mathcal{L}I$
$\frac{1}{2}mv^2$ (kinetic energy)	$\frac{1}{2}\mathcal{L}I^2$ (magnetic energy)

17-8 Inductance and magnetic energy

Continuing with the analogy of the preceding section, we would expect that corresponding to the mechanical momentum $p = mv$, whose rate of change is the applied force, there should be an analogous quantity equal to $\mathcal{L}I$, whose rate of change is \mathcal{V} . We have no right, of course, to say that $\mathcal{L}I$ is the real momentum of the circuit; in fact, it isn't. The whole circuit may be standing still and have no momentum. It is only that $\mathcal{L}I$ is analogous to the momentum mv in the sense of satisfying corresponding equations. In the same way, to the kinetic energy $\frac{1}{2}mv^2$, there corresponds an analogous quantity $\frac{1}{2}\mathcal{L}I^2$. But there we have a surprise. This $\frac{1}{2}\mathcal{L}I^2$ is really the energy in the electrical case also. This is because the rate of doing work on the inductance is $\mathcal{V}I$, and in the mechanical system it is Fv , the corresponding quantity. Therefore, in the case of the energy, the quantities not only correspond mathematically, but also have the same physical meaning as well.

We may see this in more detail as follows. As we found in Eq. (17.16), the rate of electrical work by induced forces is the product of the electromotive force and the current:

$$\frac{dW}{dt} = \mathcal{E}I.$$

Replacing \mathcal{E} by its expression in terms of the current from Eq. (17.34), we have

$$\frac{dW}{dt} = -\mathcal{L}I \frac{dI}{dt}. \quad (17.36)$$

Integrating this equation, we find that the energy required from an external source to overcome the emf in the self-inductance while building up the current† (which must equal the energy stored, U) is

$$-W = U = \frac{1}{2}\mathcal{L}I^2 \quad (17.37)$$

Therefore the energy stored in an inductance is $\frac{1}{2}\mathcal{L}I^2$.

* This is, incidentally, *not* the *only* way a correspondence can be set up between mechanical and electrical quantities.

† We are neglecting any energy loss to heat from the current in the resistance of the coil. Such losses require additional energy from the source but do not change the energy which goes into the inductance.

Applying the same arguments to a pair of coils such as those in Figs. 17-8 or 17-9, we can show that the total electrical energy of the system is given by

$$U = \frac{1}{2}\mathcal{L}_1 I_1^2 + \frac{1}{2}\mathcal{L}_2 I_2^2 + \mathfrak{M} I_1 I_2. \quad (17.38)$$

For, starting with $I = 0$ in both coils, we could first turn on the current I_1 in coil 1, with $I_2 = 0$. The work done is just $\frac{1}{2}\mathcal{L}_1 I_1^2$. But now, on turning up I_2 , we not only do the work $\frac{1}{2}\mathcal{L}_2 I_2^2$ against the emf in circuit 2, but also an additional amount $\mathfrak{M} I_1 I_2$, which is the integral of the emf $[\mathfrak{M}(dI_2/dt)]$ in circuit 1 times the now *constant* current I_1 in that circuit.

Suppose we now wish to find the force between any two coils carrying the currents I_1 and I_2 . We might at first expect that we could use the principle of virtual work, by taking the change in the energy of Eq. (17.38). We must remember, of course, that as we change the relative positions of the coils the only quantity which varies is the mutual inductance \mathfrak{M} . We might then write the equation of virtual work as

$$-F \Delta x = \Delta U = I_1 I_2 \Delta \mathfrak{M} \text{ (wrong).}$$

But this equation is wrong because, as we have seen earlier, it includes only the change in the energy of the two coils and not the change in the energy of the sources which are maintaining the currents I_1 and I_2 at their constant values. We can now understand that these sources must supply energy against the induced emf's in the coils as they are moved. If we wish to apply the principle of virtual work correctly, we must also include these energies. As we have seen, however, we may take a short cut and use the principle of virtual work by remembering that the total energy is the negative of what we have called U_{mech} , the "mechanical energy." We can therefore write for the force

$$-F \Delta x = \Delta U_{\text{mech}} = -\Delta U. \quad (17.39)$$

The force between two coils is then given by

$$F \Delta x = I_1 I_2 \Delta \mathfrak{M}.$$

Equation (17.38) for the energy of a system of two coils can be used to show that an interesting inequality exists between mutual inductance \mathfrak{M} and the self-inductances \mathcal{L}_1 and \mathcal{L}_2 of the two coils. It is clear that the energy of two coils must be positive. If we begin with zero currents in the coils and increase these currents to some values, we have been adding energy to the system. If not, the currents would spontaneously increase with release of energy to the rest of the world—an unlikely thing to happen! Now our energy equation, Eq. (17.38), can equally well be written in the following form:

$$U = \frac{1}{2} \mathcal{L}_1 \left(I_1 + \frac{\mathfrak{M}}{\mathcal{L}_1} I_2 \right)^2 + \frac{1}{2} \left(\mathcal{L}_2 - \frac{\mathfrak{M}^2}{\mathcal{L}_1} \right) I_2^2. \quad (17.40)$$

That is just an algebraic transformation. This quantity must always be positive for any values of I_1 and I_2 . In particular, it must be positive if I_2 should happen to have the special value

$$I_2 = -\frac{\mathcal{L}_1}{\mathfrak{M}} I_1 \quad (17.41)$$

But with this current for I_2 , the first term in Eq. (17.40) is zero. If the energy is to be positive, the last term in (17.40) must be greater than zero. We have the requirement that

$$\mathcal{L}_1 \mathcal{L}_2 > \mathfrak{M}^2.$$

We have thus proved the general result that the magnitude of the mutual inductance \mathfrak{M} of any two coils is necessarily less than or equal to the geometric mean of the two self-inductances. (\mathfrak{M} itself may be positive or negative, depending on the sign

conventions for the currents I_1 and I_2 .)

$$|\mathfrak{M}| < \sqrt{\mathcal{L}_1 \mathcal{L}_2}. \quad (17.42)$$

The relation between \mathfrak{M} and the self-inductances is usually written as

$$\mathfrak{M} = k\sqrt{\mathcal{L}_1 \mathcal{L}_2}. \quad (17.43)$$

The constant k is called the coefficient of coupling. If most of the flux from one coil links the other coil, the coefficient of coupling is near one; we say the coils are “tightly coupled.” If the coils are far apart or otherwise arranged so that there is very little mutual flux linkage, the coefficient of coupling is near zero and the mutual inductance is very small.

For calculating the mutual inductance of two coils, we have given in Eq. (17.30) a formula which is a double line integral around the two circuits. We might think that the same formula could be used to get the self-inductance of a single coil by carrying out both line integrals around the same coil. This, however, will not work, because in integrating around the two coils, the denominator r_{12} of the integrand will go to zero when the two line elements are at the same point. The self-inductance obtained from this formula is infinite. The reason is that this formula is an approximation that is valid only when the cross sections of the wires of the two circuits are small compared with the distance from one circuit to the other. Clearly, this approximation doesn’t hold for a single coil. It is, in fact, true that the inductance of a single coil tends logarithmically to infinity as the diameter of its wire is made smaller and smaller.

We must, then, look for a different way of calculating the self-inductance of a single coil. It is necessary to take into account the distribution of the currents within the wires because the size of the wire is an important parameter. We should therefore ask not what is the inductance of a “circuit,” but what is the inductance of a *distribution* of conductors. Perhaps the easiest way to find this inductance is to make use of the magnetic energy. We found earlier, in Section 15-3, an expression for the magnetic energy of a distribution of stationary currents:

$$U = \frac{1}{2} \int \mathbf{j} \cdot \mathbf{A} \, dV. \quad (17.44)$$

If we know the distribution of current density \mathbf{j} , we can compute the vector potential \mathbf{A} and then evaluate the integral of Eq. (17.44) to get the energy. This energy is equal to the magnetic energy of the self-inductance, $\frac{1}{2}\mathcal{L}I^2$. Equating the two gives us a formula for the inductance:

$$\mathcal{L} = \frac{1}{I^2} \int \mathbf{j} \cdot \mathbf{A} \, dV. \quad (17.45)$$

We expect, of course, that the inductance is a number depending only on the geometry of the circuit and not on the current I in the circuit. The formula of Eq. (17.45) will indeed give such a result, because the integral in this equation is proportional to the square of the current—the current appears once through \mathbf{j} and again through the vector potential \mathbf{A} . The integral divided by I^2 will depend on the geometry of the circuit but not on the current I .

Equation (17.44) for the energy of a current distribution can be put in a quite different form which is sometimes more convenient for calculation. Also, as we will see later, it is a form that is important because it is more generally valid. In the energy equation, Eq. (17.44), both \mathbf{A} and \mathbf{j} can be related to \mathbf{B} , so we can hope to express the energy in terms of the magnetic field—just as we were able to relate the electrostatic energy to the electric field. We begin by replacing \mathbf{j} by $\epsilon_0 c^2 \nabla \times \mathbf{B}$. We cannot replace \mathbf{A} so easily, since $\mathbf{B} = \nabla \times \mathbf{A}$ cannot be reversed to give \mathbf{A} in terms of \mathbf{B} . Anyway, we can write

$$U = \frac{\epsilon_0 c^2}{2} \int (\nabla \times \mathbf{B}) \cdot \mathbf{A} \, dV. \quad (17.46)$$

The interesting thing is that—with some restrictions—this integral can be written as

$$U = \frac{\epsilon_0 c^2}{2} \int \mathbf{B} \cdot (\nabla \times \mathbf{A}) dV. \quad (17.47)$$

To see this, we write out in detail a typical term. Suppose that we take the term $(\nabla \times \mathbf{B})_z A_z$ which occurs in the integral of Eq. (17.46). Writing out the components, we get

$$\int \left(\frac{\partial B_y}{\partial x} - \frac{\partial B_x}{\partial y} \right) A_z dx dy dz.$$

(There are, of course, two more integrals of the same kind.) We now integrate the first term with respect to x —integrating by parts. That is, we can say

$$\int \frac{\partial B_y}{\partial x} A_z dx = B_y A_z - \int B_y \frac{\partial A_z}{\partial x} dx.$$

Now suppose that our system—meaning the sources and fields—is finite, so that as we go to large distances all fields go to zero. Then if the integrals are carried out over all space, evaluating the term $B_y A_z$ at the limits will give zero. We have left only the term with $B_y (\partial A_z / \partial x)$, which is evidently one part of $B_y (\nabla \times \mathbf{A})_y$ and, therefore, of $\mathbf{B} \cdot (\nabla \times \mathbf{A})$. If you work out the other five terms, you will see that Eq. (17.47) is indeed equivalent to Eq. (17.46).

But now we can replace $(\nabla \times \mathbf{A})$ by \mathbf{B} , to get

$$U = \frac{\epsilon_0 c^2}{2} \int \mathbf{B} \cdot \mathbf{B} dV. \quad (17.48)$$

We have expressed the energy of a magnetostatic situation in terms of the magnetic field only. The expression corresponds closely to the formula we found for the electrostatic energy:

$$U = \frac{\epsilon_0}{2} \int \mathbf{E} \cdot \mathbf{E} dV. \quad (17.49)$$

One reason for emphasizing these two energy formulas is that sometimes they are more convenient to use. More important, it turns out that for dynamic fields (when \mathbf{E} and \mathbf{B} are changing with time) the two expressions (17.48) and (17.49) remain true, whereas the other formulas we have given for electric or magnetic energies are no longer correct—they hold only for static fields.

If we know the magnetic field \mathbf{B} of a single coil, we can find the self-inductance by equating the energy expression (17.48) to $\frac{1}{2} \mathcal{L} I^2$. Let's see how this works by finding the self-inductance of a long solenoid. We have seen earlier that the magnetic field inside a solenoid is uniform and \mathbf{B} outside is zero. The magnitude of the field inside is $B = nI/\epsilon_0 c^2$, where n is the number of turns per unit length in the winding and I is the current. If the radius of the coil is r and its length is L (we take L very long, so that we can neglect end effects, i.e., $L \gg r$), the volume inside is $\pi r^2 L$. The magnetic energy is therefore

$$U = \frac{\epsilon_0 c^2}{2} B^2 \cdot (\text{Vol}) = \frac{n^2 I^2}{2 \epsilon_0 c^2} \pi r^2 L,$$

which is equal to $\frac{1}{2} \mathcal{L} I^2$. Or,

$$\mathcal{L} = \frac{\pi r^2 n^2}{\epsilon_0 c^2} L. \quad (17.50)$$

The Maxwell Equations

18-1 Maxwell's equations

In this chapter we come back to the complete set of the four Maxwell equations that we took as our starting point in Chapter 1. Until now, we have been studying Maxwell's equations in bits and pieces; it is time to add one final piece, and to put them all together. We will then have the complete and correct story for electromagnetic fields that may be changing with time in any way. Anything said in this chapter that contradicts something said earlier is true and what was said earlier is false—because what was said earlier applied to such special situations as, for instance, steady currents or fixed charges. Although we have been very careful to point out the restrictions whenever we wrote an equation, it is easy to forget all of the qualifications and to learn too well the wrong equations. Now we are ready to give the whole truth, with no qualifications (or almost none).

The complete Maxwell equations are written in Table 18-1, in words as well as in mathematical symbols. The fact that the words are equivalent to the equations should by this time be familiar—you should be able to translate back and forth from one form to the other.

The first equation—that the divergence of \mathbf{E} is the charge density over ϵ_0 —is true in general. In dynamic as well as in static fields, Gauss' law is always valid. The flux of \mathbf{E} through any closed surface is proportional to the charge inside. The third equation is the corresponding general law for magnetic fields. Since there are no magnetic charges, the flux of \mathbf{B} through any closed surface is always zero. The second equation, that the curl of \mathbf{E} is $-\partial\mathbf{B}/\partial t$, is Faraday's law and was discussed in the last two chapters. It also is generally true. The last equation has something new. We have seen before only the part of it which holds for steady currents. In that case we said that the curl of \mathbf{B} is $\mathbf{j}/\epsilon_0 c^2$, but the correct general equation has a new part that was discovered by Maxwell.

Until Maxwell's work, the known laws of electricity and magnetism were those we have studied in Chapters 3 through 17. In particular, the equation for the magnetic field of steady currents was known only as

$$\nabla \times \mathbf{B} = \frac{\mathbf{j}}{\epsilon_0 c^2}. \quad (18.1)$$

Maxwell began by considering these known laws and expressing them as differential equations, as we have done here. (Although the ∇ notation was not yet invented, it is mainly due to Maxwell that the importance of the combinations of derivatives, which we today call the curl and the divergence, first became apparent.) He then noticed that there was something strange about Eq. (18.1). If one takes the divergence of this equation, the left-hand side will be zero, because the divergence of a curl is always zero. So this equation requires that the divergence of \mathbf{j} also be zero. But if the divergence of \mathbf{j} is zero, then the total flux of current out of any closed surface is also zero.

The flux of current from a closed surface is the decrease of the charge inside the surface. This certainly cannot in general be zero because we know that the charges can be moved from one place to another. The equation

$$\nabla \cdot \mathbf{j} = -\frac{\partial \rho}{\partial t} \quad (18.2)$$

has, in fact, been almost our definition of \mathbf{j} . This equation expresses the very funda-

18-1 Maxwell's equations

18-2 How the new term works

18-3 All of classical physics

18-4 A travelling field

18-5 The speed of light

18-6 Solving Maxwell's equations; the potentials and the wave equation

Table 18-1 Classical Physics

Maxwell's equations

- I. $\nabla \cdot \mathbf{E} = \frac{\rho}{\epsilon_0}$ (Flux of \mathbf{E} through a closed surface) = (Charge inside)/ ϵ_0
- II. $\nabla \times \mathbf{E} = -\frac{\partial \mathbf{B}}{\partial t}$ (Line integral of \mathbf{E} around a loop) = $-\frac{d}{dt}$ (Flux of \mathbf{B} through the loop)
- III. $\nabla \cdot \mathbf{B} = 0$ (Flux of \mathbf{B} through a closed surface) = 0
- IV. $c^2 \nabla \times \mathbf{B} = \frac{\mathbf{j}}{\epsilon_0} + \frac{\partial \mathbf{E}}{\partial t}$ c^2 (Integral of \mathbf{B} around a loop) = (Current through the loop)/ ϵ_0
 $+ \frac{\partial}{\partial t}$ (Flux of \mathbf{E} through the loop)

$$\left[\begin{array}{l} \text{Conservation of charge} \\ \nabla \cdot \mathbf{j} = -\frac{\partial \rho}{\partial t} \end{array} \quad \begin{array}{l} \text{(Flux of current through a closed surface)} = -\frac{\partial}{\partial t} \text{ (Charge inside)} \end{array} \right]$$

Force law

$$\mathbf{F} = q(\mathbf{E} + \mathbf{v} \times \mathbf{B})$$

Law of motion

$$\frac{d}{dt}(\mathbf{p}) = \mathbf{F}, \quad \text{where} \quad \mathbf{p} = \frac{m\mathbf{v}}{\sqrt{1 - v^2/c^2}} \quad \text{(Newton's law, with Einstein's modification)}$$

Gravitation

$$\mathbf{F} = -G \frac{m_1 m_2}{r^2} \mathbf{e}_r$$

mental law that electric charge is conserved—any flow of charge must come from some supply. Maxwell appreciated this difficulty and proposed that it could be avoided by adding the term $\partial \mathbf{E}/\partial t$ to the right-hand side of Eq. (18.1); he then got the fourth equation in Table 18-1:

$$\text{IV.} \quad c^2 \nabla \times \mathbf{B} = \frac{\mathbf{j}}{\epsilon_0} + \frac{\partial \mathbf{E}}{\partial t}.$$

It was not yet customary in Maxwell's time to think in terms of abstract fields. Maxwell discussed his ideas in terms of a model in which the vacuum was like an elastic solid. He also tried to explain the meaning of his new equation in terms of the mechanical model. There was much reluctance to accept his theory, first because of the model, and second because there was at first no experimental justification. Today, we understand better that what counts are the equations themselves and not the model used to get them. We may only question whether the equations are true or false. This is answered by doing experiments, and untold numbers of experiments have confirmed Maxwell's equations. If we take away the scaffolding he used to build it, we find that Maxwell's beautiful edifice stands on its own. He brought together all of the laws of electricity and magnetism and made one complete and beautiful theory.

Let us show that the extra term is just what is required to straighten out the difficulty Maxwell discovered. Taking the divergence of his equation (IV in Table 18-1), we must have that the divergence of the right-hand side is zero:

$$\nabla \cdot \frac{\mathbf{j}}{\epsilon_0} + \nabla \cdot \frac{\partial \mathbf{E}}{\partial t} = 0. \quad (18.3)$$

In the second term, the order of the derivatives with respect to coordinates and time can be reversed, so the equation can be rewritten as

$$\nabla \cdot \mathbf{j} + \epsilon_0 \frac{\partial}{\partial t} \nabla \cdot \mathbf{E} = 0. \quad (18.4)$$

But the first of Maxwell's equations says that the divergence of \mathbf{E} is ρ/ϵ_0 . Inserting this equality in Eq. (18.4), we get back Eq. (18.2), which we know is true. Conversely, if we accept Maxwell's equations—and we do because no one has ever found an experiment that disagrees with them—we must conclude that charge is always conserved.

The laws of physics have no answer to the question: “What happens if a charge is suddenly created at this point—what electromagnetic effects are produced?” No answer can be given because our equations say it doesn't happen. If it *were* to happen, we would need new laws, but we cannot say what they would be. We have not had the chance to observe how a world without charge conservation behaves. According to our equations, if you suddenly place a charge at some point, you had to carry it there from somewhere else. In that case, we can say what would happen.

When we added a new term to the equation for the curl of \mathbf{E} , we found that a whole new class of phenomena was described. We shall see that Maxwell's little addition to the equation for $\nabla \times \mathbf{B}$ also has far-reaching consequences. We can touch on only a few of them in this chapter.

18-2 How the new term works

As our first example we consider what happens with a spherically symmetric radial distribution of current. Suppose we imagine a little sphere with radioactive material on it. This radioactive material is squirting out some charged particles. (Or we could imagine a large block of jello with a small hole in the center into which some charge had been injected with a hypodermic needle and from which the charge is slowly leaking out.) In either case we would have a current that is everywhere radially outward. We will assume that it has the same magnitude in all directions.

Let the total charge inside any radius r be $Q(r)$. If the radial current density at the same radius is $\mathbf{j}(r)$, then Eq. (18.2) requires that Q decreases at the rate

$$\frac{\partial Q(r)}{\partial t} = -4\pi r^2 j(r). \quad (18.5)$$

We now ask about the magnetic field produced by the currents in this situation. Suppose we draw some loop Γ on a sphere of radius r , as shown in Fig. 18-1. There is some current through this loop, so we might expect to find a magnetic field circulating in the direction shown.

But we are already in difficulty. How can the \mathbf{B} have any particular direction on the sphere? A different choice of Γ would allow us to conclude that its direction is exactly opposite to that shown. So how *can* there be any circulation of \mathbf{B} around the currents?

We are saved by Maxwell's equation. The circulation of \mathbf{B} depends not only on the total *current* through Γ but also on the rate of change with time of the *electric flux* through it. It must be that these two parts just cancel. Let's see if that works out.

The electric field at the radius r must be $Q(r)/4\pi\epsilon_0 r^2$ —so long as the charge is symmetrically distributed, as we assume. It is radial, and its rate of change is then

$$\frac{\partial \mathbf{E}}{\partial t} = \frac{1}{4\pi\epsilon_0 r^2} \frac{\partial Q}{\partial t}. \quad (18.6)$$

Comparing this with Eq. (18.5), we see that at any radius

$$\frac{\partial \mathbf{E}}{\partial t} = - \frac{\mathbf{j}}{\epsilon_0}. \quad (18.7)$$

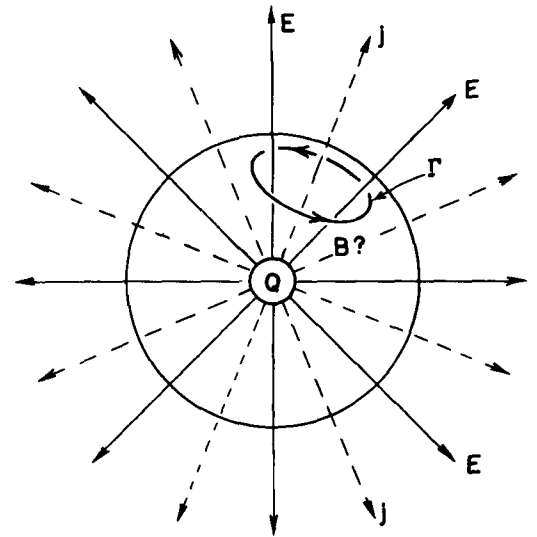


Fig. 18-1. What is the magnetic field of a spherically symmetric current?

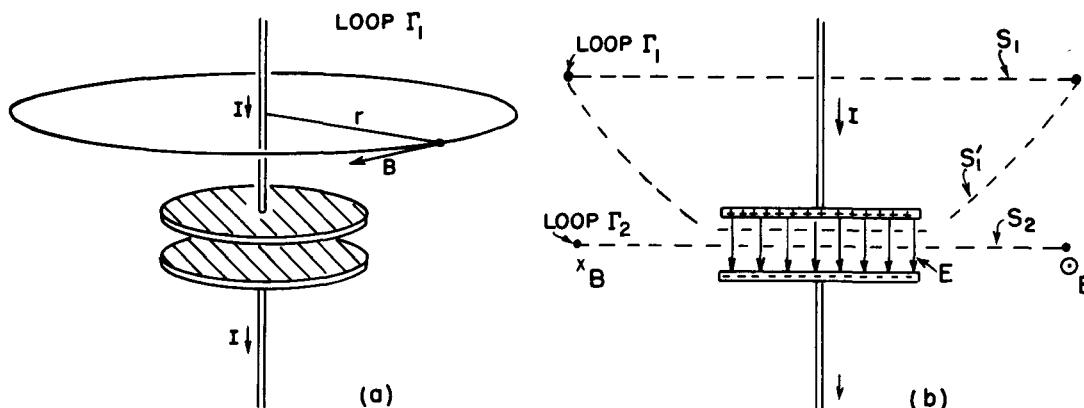


Fig. 18-2. The magnetic field near a charging capacitor.

In Eq. IV the two source terms cancel and the curl of \mathbf{B} is always zero. There is no magnetic field in our example.

As our second example, we consider the magnetic field of a wire used to charge a parallel-plate condenser (see Fig. 18-2). If the charge Q on the plates is changing with time (but not too fast), the current in the wires is equal to dQ/dt . We would expect that this current will produce a magnetic field that encircles the wire. Surely, the current close to the wire must produce the normal magnetic field—it cannot depend on where the current is going.

Suppose we take a loop Γ_1 which is a circle with radius r , as shown in part (a) of the figure. The line integral of the magnetic field should be equal to the current I divided by $\epsilon_0 c^2$. We have

$$2\pi r B = \frac{I}{\epsilon_0 c^2}. \quad (18.8)$$

This is what we would get for a steady current, but it is also correct with Maxwell's addition, because if we consider the plane surface S inside the circle, there are no electric fields on it (assuming the wire to be a very good conductor). The surface integral of $\partial \mathbf{E}/\partial t$ is zero.

Suppose, however, that we now slowly move the curve Γ downward. We get always the same result until we draw even with the plates of the condenser. Then the current I goes to zero. Does the magnetic field disappear? That would be quite strange. Let's see what Maxwell's equation says for the curve Γ_2 , which is a circle of radius r whose plane passes between the condenser plates [Fig. 18-2(b)]. The line integral of \mathbf{B} around Γ_2 is $2\pi r B$. This must equal the time derivative of the flux of \mathbf{E} through the plane circular surface S_2 . This flux of \mathbf{E} , we know from Gauss' law, must be equal to $1/\epsilon_0$ times the charge Q on one of the condenser plates. We have

$$c^2 2\pi r B = \frac{d}{dt} \left(\frac{Q}{\epsilon_0} \right). \quad (18.9)$$

That is very convenient. It is the same result we found in Eq. (18.8). Integrating over the changing electric field gives the same magnetic field as does integrating over the current in the wire. Of course, that is just what Maxwell's equation says. It is easy to see that this must always be so by applying our same arguments to the two surfaces S_1 and S'_1 that are bounded by the same circle Γ_1 in Fig. 18-2(b). Through S_1 there is the current I , but no electric flux. Through S'_1 there is no current, but an electric flux changing at the rate I/ϵ_0 . The same \mathbf{B} is obtained if we use Eq. IV with either surface.

From our discussion so far of Maxwell's new term, you may have the impression that it doesn't add much—that it just fixes up the equations to agree with what we already expect. It is true that if we just consider Eq. IV *by itself*, nothing particularly new comes out. The words "*by itself*" are, however, all-important. Maxwell's small change in Eq. IV, when *combined with the other* equations, does indeed produce much that is new and important. Before we take up these matters, however, we want to speak more about Table 18-1.

18-3 All of classical physics

In Table 18-1 we have all that was known of fundamental *classical* physics, that is, the physics that was known by 1905. Here it all is, in one table. With these equations we can understand the complete realm of classical physics.

First we have the Maxwell equations—written in both the expanded form and the short mathematical form. Then there is the conservation of charge, which is even written in parentheses, because the moment we have the complete Maxwell equations, we can deduce from them the conservation of charge. So the table is even a little redundant. Next, we have written the force law, because having all the electric and magnetic fields doesn't tell us anything until we know what they do to charges. Knowing E and B , however, we can find the force on an object with the charge q moving with velocity v . Finally, having the force doesn't tell us anything until we know what happens when a force pushes on something; we need the law of motion, which is that the force is equal to the rate of change of the momentum. (Remember? We had that in Volume I.) We even include relativity effects by writing the momentum as $p = m_0 v / \sqrt{1 - v^2/c^2}$.

If we really want to be complete, we should add one more law—Newton's law of gravitation—so we put that at the end.

Therefore in one small table we have all the fundamental laws of classical physics—even with room to write them out in words and with some redundancy. This is a great moment. We have climbed a great peak. We are on the top of K-2—we are nearly ready for Mount Everest, which is quantum mechanics. We have climbed the peak of a "Great Divide," and now we can go down the other side.

We have mainly been trying to learn how to understand the equations. Now that we have the whole thing put together, we are going to study what the equations mean—what new things they say that we haven't already seen. We've been working hard to get up to this point. It has been a great effort, but now we are going to have nice coasting downhill as we see all the consequences of our accomplishment.

18-4 A travelling field

Now for the new consequences. They come from putting together all of Maxwell's equations. First, let's see what would happen in a circumstance which we pick to be particularly simple. By assuming that all the quantities vary only in one coordinate, we will have a one-dimensional problem. The situation is shown in Fig. 18-3. We have a sheet of charge located on the yz -plane. The sheet is first at rest, then instantaneously given a velocity u in the y -direction, and kept moving with this constant velocity. You might worry about having such an "infinite" acceleration, but it doesn't really matter; just imagine that the velocity is brought to u very quickly. So we have suddenly a surface current J (J is the current per unit

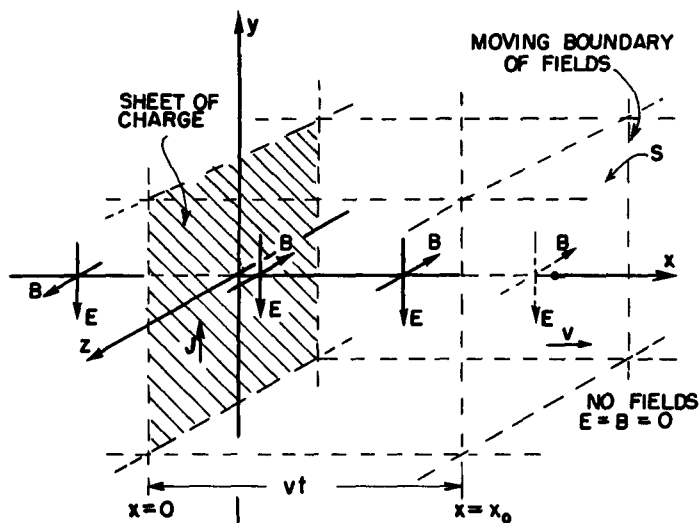


Fig. 18-3. An infinite sheet of charge is suddenly set into motion parallel to itself. There are magnetic and electric fields that propagate out from the sheet at a constant speed.

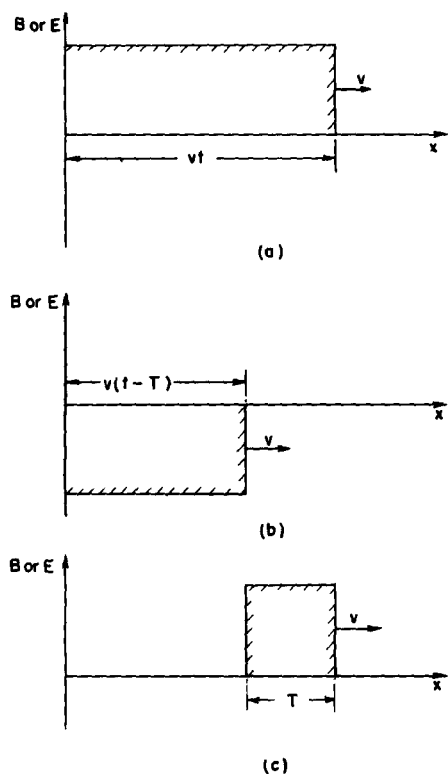


Fig. 18-4. (a) The magnitude of B (or E) as a function of x at the time t after the charge sheet is set in motion. (b) The fields for a charge sheet set in motion, toward negative y at $t = T$. (c) The sum of (a) and (b).

width in the z -direction). To keep the problem simple, we suppose that there is also a stationary sheet of charge of opposite sign superposed on the yz -plane, so that there are no electrostatic effects. Also, although in the figure we show only what is happening in a finite region, we imagine that the sheet extends to infinity in $\pm y$ and $\pm z$. In other words, we have a situation where there is no current, and then suddenly there is a uniform sheet of current. What will happen?

Well, when there is a sheet of current in the plus y -direction, there is, as we know, a magnetic field generated which will be in the minus z -direction for $x > 0$ and in the opposite direction for $x < 0$. We could find the magnitude of B by using the fact that the line integral of the magnetic field will be equal to the current over $\epsilon_0 c^2$. We would get that $B = J/2\epsilon_0 c^2$ (since the current I in a strip of width w is Jw and the line integral of B is $2Bw$).

This gives us the field next to the sheet—for small x —but since we are imagining an infinite sheet, we would expect the same argument to give the magnetic field farther out for larger values of x . However, that would mean that the moment we turn on the current, the magnetic field is suddenly changed from zero to a finite value everywhere. But wait! If the magnetic field is suddenly changed, it will produce tremendous electrical effects. (If it changes in *any* way, there are electrical effects.) So because we moved the sheet of charge, we make a changing magnetic field, and therefore electric fields must be generated. If there are electric fields generated, they had to start from zero and change to something else. There will be some $\partial E/\partial t$ that will make a contribution, together with the current J , to the production of the magnetic field. So through the various equations there is a big intermixing, and we have to try to solve for all the fields at once.

By looking at the Maxwell equations alone, it is not easy to see directly how to get the solution. So we will first show you what the answer is and then verify that it does indeed satisfy the equations. The answer is the following: The field B that we computed is, in fact, generated right next to the current sheet (for small x). It must be so, because if we make a tiny loop around the sheet, there is no room for any electric flux to go through it. But the field B out farther—for larger x —is, at first, zero. It stays zero for awhile, and then suddenly turns on. In short, we turn on the current and the magnetic field immediately next to it turns on to a constant value B ; then the turning on of B spreads out from the source region. After a certain time, there is a uniform magnetic field everywhere out to some value x , and then zero beyond. Because of the symmetry, it spreads in both the plus and minus x -directions.

The E -field does the same thing. Before $t = 0$ (when we turn on the current), the field is zero everywhere. Then after the time t , both E and B are uniform out to the distance $x = vt$, and zero beyond. The fields make their way forward like a tidal wave, with a front moving at a uniform velocity which turns out to be c , but for a while we will just call it v . A graph of the magnitude of E or B versus x , as they appear at the time t , is shown in Fig. 18-4(a). Looking again at Fig. 18-3, at the time t , the region between $x = \pm vt$ is “filled” with the fields, but they have not yet reached beyond. We emphasize again that we are assuming that the current sheet and, therefore the fields E and B , extend infinitely far in both the y - and z -directions. (We cannot draw an infinite sheet, so we have shown only what happens in a finite area.)

We want now to analyze quantitatively what is happening. To do that, we want to look at two cross-sectional views, a top view looking down along the y -axis, as shown in Fig. 18-5, and a side view looking back along the z -axis, as shown in Fig. 18-6. Suppose we start with the side view. We see the charged sheet moving up; the magnetic field points into the page for $+x$, and out of the page for $-x$, and the electric field is downward everywhere—out to $x = \pm vt$.

Let's see if these fields are consistent with Maxwell's equations. Let's first draw one of those loops that we use to calculate a line integral, say the rectangle Γ_2 shown in Fig. 18-6. You notice that one side of the rectangle is in the region where there are fields, but one side is in the region the fields have still not reached. There is some magnetic flux through this loop. If it is changing, there should be an emf around it. If the wavefront is moving, we will have a changing magnetic

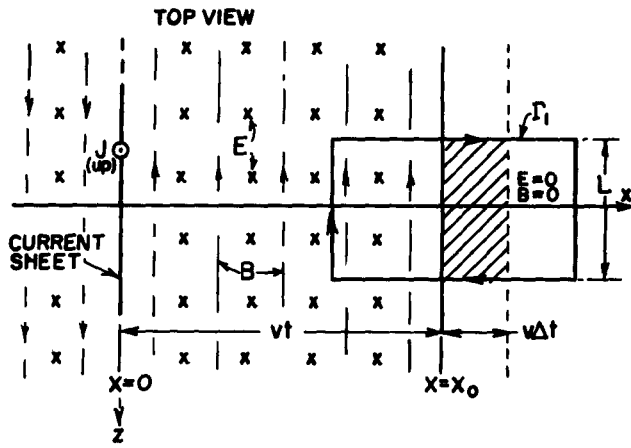


Fig. 18-5. Top view of Fig. 18-3.

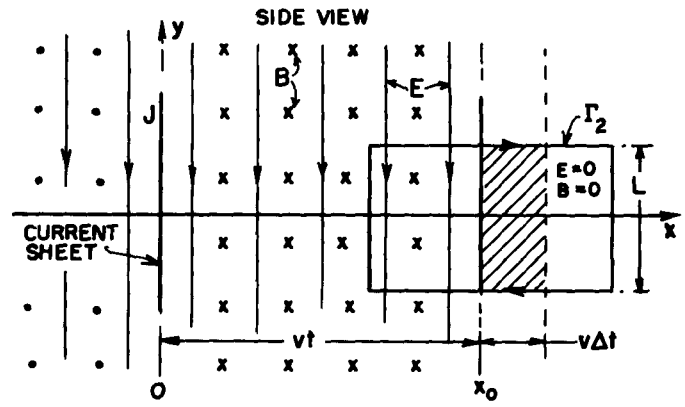


Fig. 18-6. Side view of Fig. 18-3.

flux, because the area in which B exists is progressively increasing at the velocity v . The flux inside Γ_2 is B times the part of the area inside Γ_2 which has a magnetic field. The rate of change of the flux, since the magnitude of B is constant, is the magnitude times the rate of change of the area. The rate of change of the area is easy. If the width of the rectangle Γ_2 is L , the area in which B exists changes by $Lv \Delta t$ in the time Δt . (See Fig. 18-6.) The rate of change of flux is then BLv . According to Faraday's law, this should equal the line integral of E around Γ_2 , which is just EL . We have the equation

$$E = vB. \quad (18.10)$$

So if the ratio of E to B is v , the fields we have assumed will satisfy Faraday's equation.

But that is not the only equation; we have the other equation relating E and B :

$$c^2 \nabla \times \mathbf{B} = \frac{\mathbf{j}}{\epsilon_0} + \frac{\partial \mathbf{E}}{\partial t}. \quad (18.11)$$

To apply this equation, we look at the top view in Fig. 18-5. We have seen that this equation will give us the value of B next to the current sheet. Also, for any loop drawn outside the sheet but behind the wavefront, there is no curl of B nor any j or changing E , so the equation is correct there. Now let's look at what happens for the curve Γ_1 that intersects the wavefront, as shown in Fig. 18-5. Here there are no currents, so Eq. (18.11) can be written—in integral form—as

$$c^2 \oint_{\Gamma_1} \mathbf{B} \cdot d\mathbf{s} = \frac{d}{dt} \int_{\text{inside } \Gamma_1} \mathbf{E} \cdot \mathbf{n} \, da. \quad (18.12)$$

The line integral of B is just B times L . The rate of change of the flux of E is due only to the advancing wavefront. The area inside Γ_1 , where E is not zero, is increasing at the rate vL . The right-hand side of Eq. (18.12) is then vLE . That equation becomes

$$c^2 B = Ev. \quad (18.13)$$

We have a solution in which we have a constant B and a constant E behind the front, both at right angles to the direction in which the front is moving and at right angles to each other. Maxwell's equations specify the ratio of E to B . From Eqs. (18.10) and (18.13),

$$E = vB, \quad \text{and} \quad E = \frac{c^2}{v} B.$$

But one moment! We have found *two different* conditions on the ratio E/B . Can such a field as we describe really exist? There is, of course, only one velocity v for which both of these equations can hold, namely $v = c$. The wavefront must travel with the velocity c . We have an example in which the electrical influence from a current propagates at a certain finite velocity c .

Now let's ask what happens if we suddenly stop the motion of the charged sheet after it has been on for a short time T . We can see what will happen by the principle of superposition. We had a current that was zero and then was suddenly turned on. We know the solution for that case. Now we are going to add another set of fields. We take another charged sheet and suddenly start it moving, in the opposite direction with the same speed, only at the time T after we started the first current. The total current of the two added together is first zero, then on for a time T , then off again—because the two currents cancel. We have a square “pulse” of current.

The new negative current produces the same fields as the positive one, only with all the signs reversed and, of course, delayed in time by T . A wavefront again travels out at the velocity c . At the time t it has reached the distance $x = \pm c(t - T)$, as shown in Fig. 18-4(b). So we have two “blocks” of field marching out at the speed c , as in parts (a) and (b) of Fig. 18-4. The combined fields are as shown in part (c) of the figure. The fields are zero for $x > ct$, they are constant (with the values we found above) between $x = c(t - T)$ and $x = ct$, and again zero for $x < c(t - T)$.

In short, we have a little piece of field—a block of thickness cT —which has left the current sheet and is travelling through space all by itself. The fields have “taken off”; they are propagating freely through space, no longer connected in any way with the source. The caterpillar has turned into a butterfly!

How can this bundle of electric and magnetic fields maintain itself? The answer is: by the combined effects of the Faraday law, $\nabla \times \mathbf{E} = -\partial \mathbf{B} / \partial t$, and the new term of Maxwell, $c^2 \nabla \times \mathbf{B} = \partial \mathbf{E} / \partial t$. They cannot help maintaining themselves. Suppose the magnetic field were to disappear. There would be a changing magnetic field which would produce an electric field. If this electric field tries to go away, the changing electric field would create a magnetic field back again. So by a perpetual interplay—by the swishing back and forth from one field to the other—they must go on forever. It is impossible for them to disappear.* They maintain themselves in a kind of a dance—one making the other, the second making the first—propagating onward through space.

18-5 The speed of light

We have a wave which leaves the material source and goes outward at the velocity c , which is the speed of light. But let's go back a moment. From a historical point of view, it wasn't known that the coefficient c in Maxwell's equations was also the speed of light propagation. There was just a constant in the equations. We have called it c from the beginning, because we knew what it would turn out to be. We didn't think it would be sensible to make you learn the formulas with a different constant and then go back to substitute c wherever it belonged. From the point of view of electricity and magnetism, however, we just start out with two constants, ϵ_0 and c^2 , that appear in the equations of electrostatics and magnetostatics:

$$\nabla \cdot \mathbf{E} = \frac{\rho}{\epsilon_0} \quad (18.14)$$

and

$$\nabla \times \mathbf{B} = \frac{\mathbf{j}}{\epsilon_0 c^2}. \quad (18.15)$$

If we take any *arbitrary* definition of a unit of charge, we can determine experimentally the constant ϵ_0 required in Eq. (18.14)—say by measuring the force between two unit charges at rest, using Coulomb's law. We must also determine experimentally the constant $\epsilon_0 c^2$ that appears in Eq. (18.15), which we can do, say, by measuring the force between two unit currents. (A unit current means one unit of charge per second.) The ratio of these two experimental constants is c^2 —just another “electromagnetic constant.”

* Well, not quite. They can be “absorbed” if they get to a region where there are charges. By which we mean that other fields can be produced somewhere which superpose on these fields and “cancel” them by destructive interference (see Chapter 31, Vol. I).

Notice now that this constant c^2 is the same no matter what we choose for our unit of charge. If we put twice as much “charge”—say twice as many proton charges—in our “unit” of charge, ϵ_0 would need to be one-fourth as large. When we pass two of these “unit” currents through two wires, there will be twice as much “charge” per second in each wire, so the force between two wires is four times larger. The constant $\epsilon_0 c^2$ must be reduced by one-fourth. But the ratio $\epsilon_0 c^2 / \epsilon_0$ is unchanged.

So just by experiments with charges and currents we find a number c^2 which turns out to be the square of the velocity of propagation of electromagnetic influences. From static measurements—by measuring the forces between two unit charges and between two unit currents—we find that $c = 3.00 \times 10^8$ meters/sec. When Maxwell first made this calculation with his equations, he said that bundles of electric and magnetic fields should be propagated at this speed. He also remarked on the mysterious coincidence that this was the same as the speed of light. “We can scarcely avoid the inference,” said Maxwell, “that light consists in the transverse undulations of the same medium which is the cause of electric and magnetic phenomena.”

Maxwell had made one of the great unifications of physics. Before his time, there was light, and there was electricity and magnetism. The latter two had been unified by the experimental work of Faraday, Oersted, and Ampere. Then, all of a sudden, light was no longer “something else,” but was only electricity and magnetism in this new form—little pieces of electric and magnetic fields which propagate through space on their own.

We have called your attention to some characteristics of this special solution, which turn out to be true, however, for *any* electromagnetic wave: that the magnetic field is perpendicular to the direction of motion of the wavefront; that the electric field is likewise perpendicular to the direction of motion of the wavefront; and that the two vectors \mathbf{E} and \mathbf{B} are perpendicular to each other. Furthermore, the magnitude of the electric field E is equal to c times the magnitude of the magnetic field B . These three facts—that the two fields are transverse to the direction of propagation, that \mathbf{B} is perpendicular to \mathbf{E} , and that $E = cB$ —are generally true for any electromagnetic wave. Our special case is a good one—it shows all the main features of electromagnetic waves.

18-6 Solving Maxwell's equations; the potentials and the wave equation

Now we would like to do something mathematical; we want to write Maxwell's equations in a simpler form. You may consider that we are complicating them, but if you will be patient a little bit, they will suddenly come out simpler. Although by this time you are thoroughly used to each of the Maxwell equations, there are many pieces that must all be put together. That's what we want to do.

We begin with $\nabla \cdot \mathbf{B} = 0$ —the simplest of the equations. We know that it implies that \mathbf{B} is the curl of something. So, if we write

$$\mathbf{B} = \nabla \times \mathbf{A}, \quad (18.16)$$

we have already solved one of Maxwell's equations. (Incidentally, you appreciate that it remains true that another vector \mathbf{A}' would be just as good if $\mathbf{A}' = \mathbf{A} + \nabla\psi$ —where ψ is any scalar field—because the curl of $\nabla\psi$ is zero, and \mathbf{B} is still the same. We have talked about that before.)

We take next the Faraday law, $\nabla \times \mathbf{E} = -\partial\mathbf{B}/\partial t$, because it doesn't involve any currents or charges. If we write \mathbf{B} as $\nabla \times \mathbf{A}$ and differentiate with respect to t , we can write Faraday's law in the form

$$\nabla \times \mathbf{E} = -\frac{\partial}{\partial t} \nabla \times \mathbf{A}.$$

Since we can differentiate either with respect to time or to space first, we can also write this equation as

$$\nabla \times \left(\mathbf{E} + \frac{\partial \mathbf{A}}{\partial t} \right) = 0. \quad (18.17)$$

We see that $E + \partial A/\partial t$ is a vector whose curl is equal to zero. Therefore that vector is the gradient of something. When we worked on electrostatics, we had $\nabla \times E = 0$, and then we decided that E itself was the gradient of something. We took it to be the gradient of $-\phi$ (the minus for technical convenience). We do the same thing for $E + \partial A/\partial t$; we set

$$E + \frac{\partial A}{\partial t} = -\nabla\phi. \quad (18.18)$$

We use the same symbol ϕ so that, in the electrostatic case where nothing changes with time and the $\partial A/\partial t$ term disappears, E will be our old $-\nabla\phi$. So Faraday's equation can be put in the form

$$E = -\nabla\phi - \frac{\partial A}{\partial t}. \quad (18.19)$$

We have solved two of Maxwell's equations already, and we have found that to describe the electromagnetic fields E and B , we need four potential functions: a scalar potential ϕ and a vector potential A , which is, of course, three functions.

Now that A determines part of E , as well as B , what happens when we change A to $A' = A + \nabla\psi$? In general, E would change if we didn't take some special precaution. We can, however, still allow A to be changed in this way without affecting the fields E and B —that is, without changing the physics—if we always change A and ϕ *together* by the rules

$$A' = A + \nabla\psi, \quad \phi' = \phi - \frac{\partial\psi}{\partial t}. \quad (18.20)$$

Then neither B nor E , obtained from Eq. (18.19), is changed.

Previously, we chose to make $\nabla \cdot A = 0$, to make the equations of statics somewhat simpler. We are not going to do that now; we are going to make a different choice. But we'll wait a bit before saying what the choice is, because later it will be clear *why* the choice is made.

Now we return to the two remaining Maxwell equations which will give us relations between the potentials and the sources ρ and j . Once we can determine A and ϕ from the currents and charges, we can always get E and B from Eqs. (18.16) and (18.19), so we will have another form of Maxwell's equations.

We begin by substituting Eq. (18.19) into $\nabla \cdot E = \rho/\epsilon_0$; we get

$$\nabla \cdot \left(-\nabla\phi - \frac{\partial A}{\partial t} \right) = \frac{\rho}{\epsilon_0},$$

which we can write also as

$$-\nabla^2\phi - \frac{\partial}{\partial t} \nabla \cdot A = \frac{\rho}{\epsilon_0}. \quad (18.21)$$

This is one equation relating ϕ and A to the sources.

Our final equation will be the most complicated. We start by rewriting the fourth Maxwell equation as

$$c^2 \nabla \times B - \frac{\partial E}{\partial t} = \frac{j}{\epsilon_0},$$

and then substitute for B and E in terms of the potentials, using Eqs. (18.16) and (18.19):

$$c^2 \nabla \times (\nabla \times A) - \frac{\partial}{\partial t} \left(-\nabla\phi - \frac{\partial A}{\partial t} \right) = \frac{j}{\epsilon_0}.$$

The first term can be rewritten using the algebraic identity: $\nabla \times (\nabla \times A) = \nabla(\nabla \cdot A) - \nabla^2 A$; we get

$$-c^2 \nabla^2 A + c^2 \nabla(\nabla \cdot A) + \frac{\partial}{\partial t} \nabla\phi + \frac{\partial^2 A}{\partial t^2} = \frac{j}{\epsilon_0}. \quad (18.22)$$

It's not very simple!

Fortunately, we can now make use of our freedom to choose arbitrarily the divergence of \mathbf{A} . What we are going to do is to use our choice to fix things so that the equations for \mathbf{A} and for ϕ are separated but have the same form. We can do this by taking*

$$\nabla \cdot \mathbf{A} = -\frac{1}{c^2} \frac{\partial \phi}{\partial t}. \quad (18.23)$$

When we do that, the two middle terms in \mathbf{A} and ϕ in Eq. (18.22) cancel, and that equation becomes much simpler:

$$\nabla^2 \mathbf{A} - \frac{1}{c^2} \frac{\partial^2 \mathbf{A}}{\partial t^2} = -\frac{\mathbf{j}}{\epsilon_0 c^2}. \quad (18.24)$$

And our equation for ϕ —Eq. (18.21)—takes on the same form:

$$\nabla^2 \phi - \frac{1}{c^2} \frac{\partial^2 \phi}{\partial t^2} = -\frac{\rho}{\epsilon_0}. \quad (18.25)$$

What a beautiful set of equations! They are beautiful, first, because they are nicely separated—with the charge density, goes ϕ ; with the current, goes \mathbf{A} . Furthermore, although the left side looks a little funny—a Laplacian together with a $(\partial/\partial t)^2$ —when we unfold it we see

$$\frac{\partial^2 \phi}{\partial x^2} + \frac{\partial^2 \phi}{\partial y^2} + \frac{\partial^2 \phi}{\partial z^2} - \frac{1}{c^2} \frac{\partial^2 \phi}{\partial t^2} = -\frac{\rho}{\epsilon_0}. \quad (18.26)$$

It has a nice symmetry in x, y, z, t —the $-1/c^2$ is necessary because, of course, time and space *are* different; they have different units.

Maxwell's equations have led us to a new kind of equation for the potentials ϕ and \mathbf{A} but to the same mathematical form for all four functions $\phi, A_x, A_y,$ and A_z . Once we learn how to solve these equations, we can get \mathbf{B} and \mathbf{E} from $\nabla \times \mathbf{A}$ and $-\nabla\phi - \partial\mathbf{A}/\partial t$. We have another form of the electromagnetic laws exactly equivalent to Maxwell's equations, and in many situations they are much simpler to handle.

We have, in fact, already solved an equation much like Eq. (18.26). When we studied sound in Chapter 47 of Vol. I, we had an equation of the form

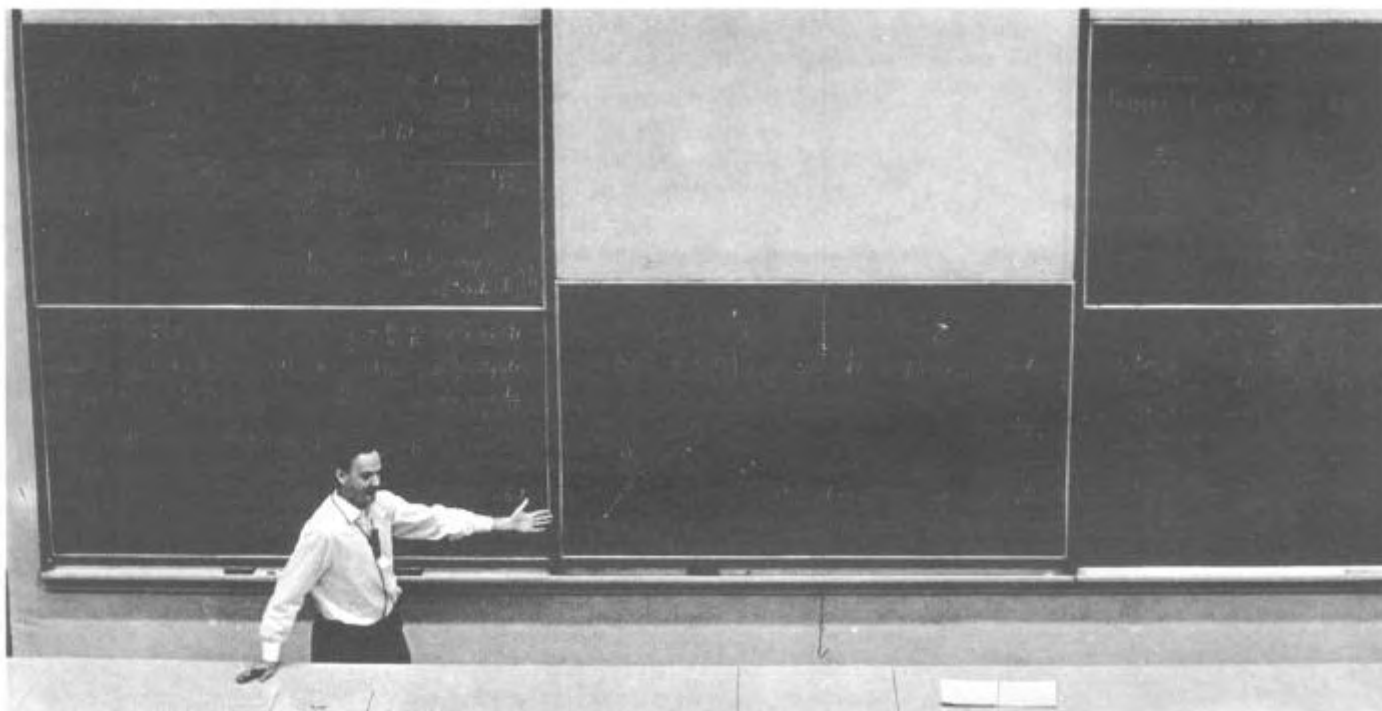
$$\frac{\partial^2 \phi}{\partial x^2} = \frac{1}{c^2} \frac{\partial^2 \phi}{\partial t^2},$$

and we saw that it described the propagation of waves in the x -direction at the speed c . Equation (18.26) is the corresponding wave equation for three dimensions. So in regions where there are no longer any charges and currents, the solution of these equations is *not* that ϕ and \mathbf{A} are zero. (Although that is indeed one possible solution.) There are solutions in which there is some set of ϕ and \mathbf{A} which are changing in time but always moving out at the speed c . The fields travel onward through free space, as in our example at the beginning of the chapter.

With Maxwell's new term in Eq. IV, we have been able to write the field equations in terms of \mathbf{A} and ϕ in a form that is simple and that makes immediately apparent that there are electromagnetic waves. For many practical purposes, it will still be convenient to use the original equations in terms of \mathbf{E} and \mathbf{B} . But they are on the other side of the mountain we have already climbed. Now we are ready to cross over to the other side of the peak. Things will look different—we are ready for some new and beautiful views.

* Choosing the $\nabla \cdot \mathbf{A}$ is called "choosing a gauge." Changing \mathbf{A} by adding $\nabla\psi$ is called a "gauge transformation." Equation (18.23) is called "the Lorentz gauge."

The Principle of Least Action



A special lecture—almost verbatim*

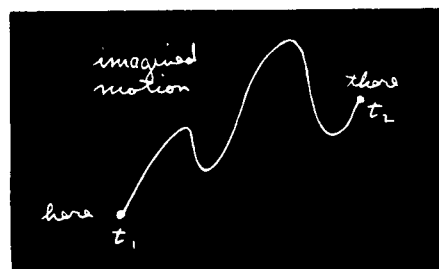
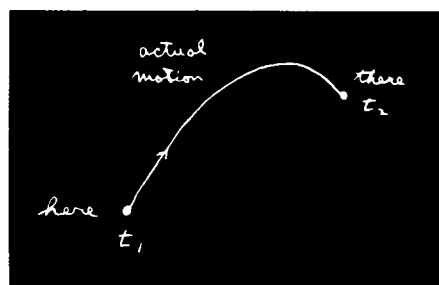
“When I was in high school, my physics teacher—whose name was Mr. Bader—called me down one day after physics class and said, ‘You look bored; I want to tell you something interesting.’ Then he told me something which I found absolutely fascinating, and have, since then, always found fascinating. Every time the subject comes up, I work on it. In fact, when I began to prepare this lecture I found myself making more analyses on the thing. Instead of worrying about the lecture, I got involved in a new problem. The subject is this—the principle of least action.

“Mr. Bader told me the following: Suppose you have a particle (in a gravitational field, for instance) which starts somewhere and moves to some other point by free motion—you throw it, and it goes up and comes down. —————→

It goes from the original place to the final place in a certain amount of time. Now, you try a different motion. Suppose that to get from here to there, it went like this —————→

but got there in just the same amount of time. Then he said this: If you calculate the kinetic energy at every moment on the path, take away the potential energy, and integrate it over the time during the whole path, you’ll find that the number you’ll get is *bigger* than that for the actual motion.

* Later chapters do not depend on the material of this special lecture—which is intended to be for “entertainment.”



"In other words, the laws of Newton could be stated not in the form $F = ma$ but in the form: the average kinetic energy less the average potential energy is as little as possible for the path of an object going from one point to another.

"Let me illustrate a little bit better what it means. If you take the case of the gravitational field, then if the particle has the path $x(t)$ (let's just take one dimension for a moment; we take a trajectory that goes up and down and not sideways), where x is the height above the ground, the kinetic energy is $\frac{1}{2}m(dx/dt)^2$, and the potential energy at any time is mgx . Now I take the kinetic energy minus the potential energy at every moment along the path and integrate that with respect to time from the initial time to the final time. Let's suppose that at the original time t_1 we started at some height and at the end of the time t_2 we are definitely ending at some other place.

"Then the integral is

$$\int_{t_1}^{t_2} \left[\frac{1}{2} m \left(\frac{dx}{dt} \right)^2 - mgx \right] dt.$$

The actual motion is some kind of a curve—it's a parabola if we plot against the time—and gives a certain value for the integral. But we could *imagine* some other motion that went very high and came up and down in some peculiar way.

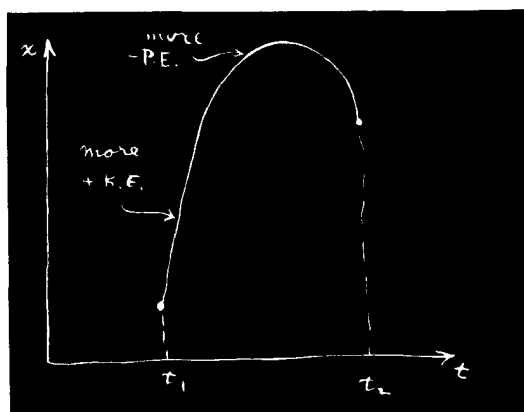
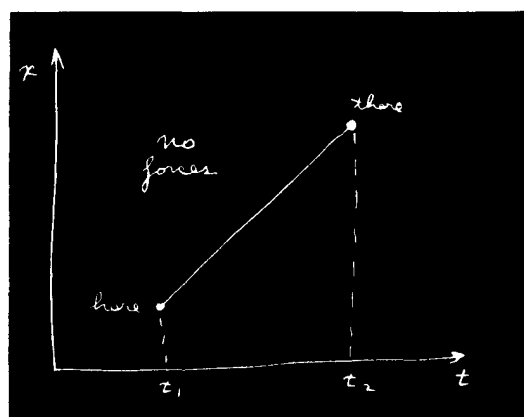
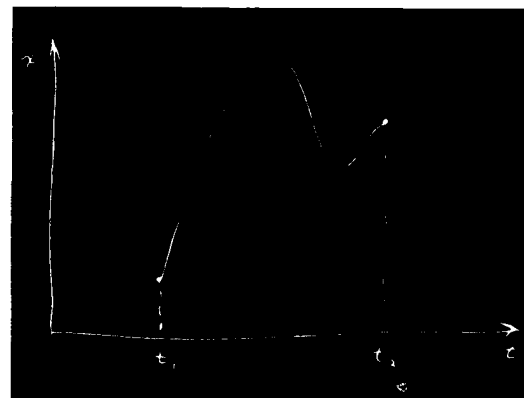
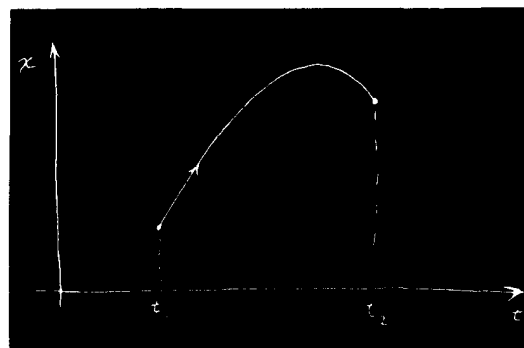
We can calculate the kinetic energy minus the potential energy and integrate for such a path . . . or for any other path we want. The miracle is that the true path is the one for which that integral is least.

"Let's try it out. First, suppose we take the case of a free particle for which there is no potential energy at all. Then the rule says that in going from one point to another in a given amount of time, the kinetic energy integral is least, so it must go at a uniform speed. (We know that's the right answer—to go at a uniform speed.) Why is that? Because if the particle were to go any other way, the velocities would be sometimes higher and sometimes lower than the average. The average velocity is the same for every case because it has to get from 'here' to 'there' in a given amount of time.

"As an example, say your job is to start from home and get to school in a given length of time with the car. You can do it several ways: You can accelerate like mad at the beginning and slow down with the brakes near the end, or you can go at a uniform speed, or you can go backwards for a while and then go forward, and so on. The thing is that the average speed has got to be, of course, the total distance that you have gone over the time. But if you do anything but go at a uniform speed, then sometimes you are going too fast and sometimes you are going too slow. Now the mean *square* of something that deviates around an average, as you know, is always greater than the square of the mean; so the kinetic energy integral would always be higher if you wobbled your velocity than if you went at a uniform velocity. So we see that the integral is a minimum if the velocity is a constant (when there are no forces). The correct path is like this.

"Now, an object thrown up in a gravitational field does rise faster first and then slow down. That is because there is also the potential energy, and we must have the least *difference* of kinetic and potential energy on the average. Because the potential energy rises as we go up in space, we will get a lower *difference* if we can get as soon as possible up to where there is a high potential energy. Then we can take that potential away from the kinetic energy and get a lower average. So it is better to take a path which goes up and gets a lot of negative stuff from the potential energy.

"On the other hand, you can't go up too fast, or too far, because you will then have too much kinetic energy involved—you have to go very fast to get way up and come down again in the fixed amount of time available. So you don't want to go too far up, but you want to go up some. So it turns out that the solution is some kind of balance between trying to get more potential energy with the least amount of extra kinetic energy—trying to get the difference, kinetic minus the potential, as small as possible.



"That is all my teacher told me, because he was a very good teacher and knew when to stop talking. But I don't know when to stop talking. So instead of leaving it as an interesting remark, I am going to horrify and disgust you with the complexities of life by proving that it is so. The kind of mathematical problem we will have is very difficult and a new kind. We have a certain quantity which is called the *action*, S . It is the kinetic energy, minus the potential energy, integrated over time.

$$\text{Action} = S = \int_{t_1}^{t_2} (\text{KE} - \text{PE}) dt.$$

Remember that the PE and KE are both functions of time. For each different possible path you get a different number for this action. Our mathematical problem is to find out for what curve that number is the least.

"You say—Oh, that's just the ordinary calculus of maxima and minima. You calculate the action and just differentiate to find the minimum.

"But watch out. Ordinarily we just have a function of some variable, and we have to find the value of that *variable* where the function is least or most. For instance, we have a rod which has been heated in the middle and the heat is spread around. For each point on the rod we have a temperature, and we must find the point at which that temperature is largest. But now for *each path in space* we have a number—quite a different thing—and we have to find the *path in space* for which the number is the minimum. That is a completely different branch of mathematics. It is not the ordinary calculus. In fact, it is called the *calculus of variations*.

"There are many problems in this kind of mathematics. For example, the circle is usually defined as the locus of all points at a constant distance from a fixed point, but another way of defining a circle is this: a circle is that curve of *given length* which encloses the biggest area. Any other curve encloses less area for a given perimeter than the circle does. So if we give the problem: find that curve which encloses the greatest area for a given perimeter, we would have a problem of the calculus of variations—a different kind of calculus than you're used to.

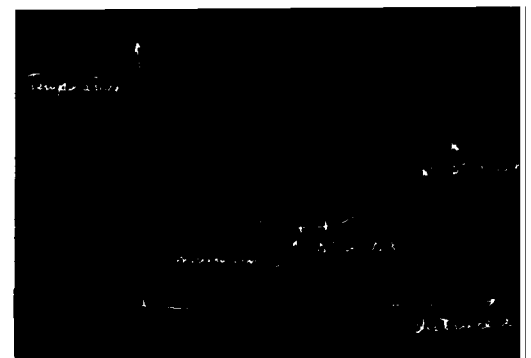
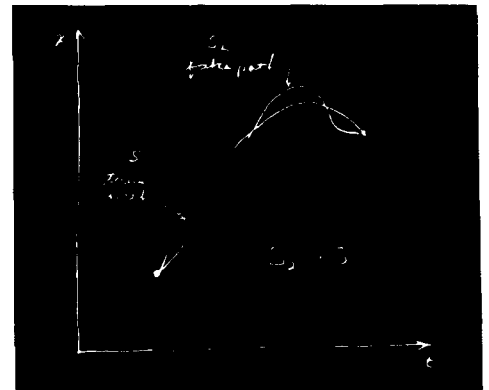
"So we make the calculation for the path of an object. Here is the way we are going to do it. The idea is that we imagine that there is a true path and that any other curve we draw is a false path, so that if we calculate the action for the false path we will get a value that is bigger than if we calculate the action for the true path.

"Problem: Find the true path. Where is it? One way, of course, is to calculate the action for millions and millions of paths and look at which one is lowest. When you find the lowest one, that's the true path.

"That's a possible way. But we can do it better than that. When we have a quantity which has a minimum—for instance, in an ordinary function like the temperature—one of the properties of the minimum is that if we go away from the minimum in the *first* order, the deviation of the function from its minimum value is only *second* order. At any place else on the curve, if we move a small distance the value of the function changes also in the first order. But at a minimum, a tiny motion away makes, in the first approximation, no difference.

"That is what we are going to use to calculate the true path. If we have the true path, a curve which differs only a little bit from it will, in the first approximation, make no difference in the action. Any difference will be in the second approximation, if we really have a minimum.

"That is easy to prove. If there is a change in the first order when I deviate the curve a certain way, there is a change in the action that is *proportional* to the deviation. The change presumably makes the action greater; otherwise we haven't got a minimum. But then if the change is *proportional* to the deviation, reversing the sign of the deviation will make the action less. We would get the action to increase one way and to decrease the other way. The only way that it could really be a minimum is that in the *first* approximation it doesn't make any change, that the changes are proportional to the square of the deviations from the true path.



“So we work it this way: We call $\underline{x(t)}$ (with an underline) the true path—the one we are trying to find. We take some trial path $x(t)$ that differs from the true path by a small amount which we will call $\eta(t)$ (eta of t).

“Now the idea is that if we calculate the action S for the path $x(t)$, then the difference between that S and the action that we calculated for the path $\underline{x(t)}$ —to simplify the writing we can call it \underline{S} —the difference of \underline{S} and S must be zero in the first-order approximation of small η . It can differ in the second order, but in the first order the difference must be zero.

“And that must be true for any η at all. Well, not quite. The method doesn't mean anything unless you consider paths which all begin and end at the same two points—each path begins at a certain point at t_1 and ends at a certain other point at t_2 , and those points and times are kept fixed. So the deviations in our η have to be zero at each end, $\eta(t_1) = 0$ and $\eta(t_2) = 0$. With that condition, we have specified our mathematical problem.

“If you didn't know any calculus, you might do the same kind of thing to find the minimum of an ordinary function $f(x)$. You could discuss what happens if you take $f(x)$ and add a small amount h to x and argue that the correction to $f(x)$ in the first order in h must be zero at the minimum. You would substitute $x + h$ for x and expand out to the first order in h . . . just as we are going to do with η .

“The idea is then that we substitute $x(t) = \underline{x(t)} + \eta(t)$ in the formula for the action:

$$S = \int \left[\frac{m}{2} \left(\frac{dx}{dt} \right)^2 - V(x) \right] dt,$$

where I call the potential energy $V(x)$. The derivative dx/dt is, of course, the derivative of $\underline{x(t)}$ plus the derivative of $\eta(t)$, so for the action I get this expression:

$$S = \int_{t_1}^{t_2} \left[\frac{m}{2} \left(\frac{d\underline{x}}{dt} + \frac{d\eta}{dt} \right)^2 - V(\underline{x} + \eta) \right] dt.$$

“Now I must write this out in more detail. For the squared term I get

$$\left(\frac{d\underline{x}}{dt} \right)^2 + 2 \frac{d\underline{x}}{dt} \frac{d\eta}{dt} + \left(\frac{d\eta}{dt} \right)^2.$$

But wait. I'm not worrying about higher than the first order, so I will take all the terms which involve η^2 and higher powers and put them in a little box called 'second and higher order.' From this term I get only second order, but there will be more from something else. So the kinetic energy part is

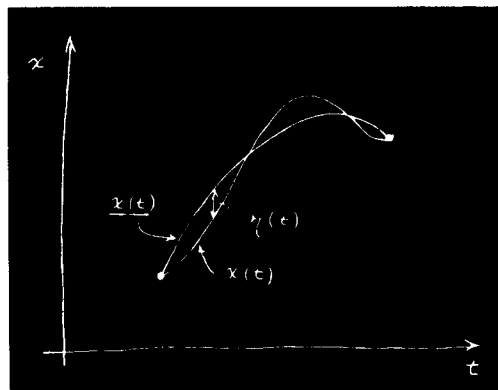
$$\frac{m}{2} \left(\frac{d\underline{x}}{dt} \right)^2 + m \frac{d\underline{x}}{dt} \frac{d\eta}{dt} + (\text{second and higher order}).$$

“Now we need the potential V at $\underline{x} + \eta$. I consider η small, so I can write $V(x)$ as a Taylor series. It is approximately $V(\underline{x})$; in the next approximation (from the ordinary nature of derivatives) the correction is η times the rate of change of V with respect to x , and so on:

$$V(\underline{x} + \eta) = V(\underline{x}) + \eta V'(\underline{x}) + \frac{\eta^2}{2} V''(\underline{x}) + \dots$$

I have written V' for the derivative of V with respect to x in order to save writing. The term in η^2 and the ones beyond fall into the 'second and higher order' category and we don't have to worry about them. Putting it all together,

$$S = \int_{t_1}^{t_2} \left[\frac{m}{2} \left(\frac{d\underline{x}}{dt} \right)^2 - V(\underline{x}) + m \frac{d\underline{x}}{dt} \frac{d\eta}{dt} - \eta V'(\underline{x}) + (\text{second and higher order}) \right] dt.$$



Now if we look carefully at the thing, we see that the first two terms which I have arranged here correspond to the action \bar{S} that I would have calculated with the true path \underline{x} . The thing I want to concentrate on is the change in S —the difference between the S and the \bar{S} that we would get for the right path. This difference we will write as δS , called the variation in S . Leaving out the ‘second and higher order’ terms, I have for δS

$$\delta S = \int_{t_1}^{t_2} \left[m \frac{dx}{dt} \frac{d\eta}{dt} - \eta V'(\underline{x}) \right] dt.$$

“Now the problem is this: Here is a certain integral. I don’t know what the \underline{x} is yet, but I do know that *no matter what* η is, this integral must be zero. Well, you think, the only way that that can happen is that what multiplies η must be zero. But what about the first term with $d\eta/dt$? Well, after all, if η can be anything at all, its derivative is anything also, so you conclude that the coefficient of $d\eta/dt$ must also be zero. That isn’t quite right. It isn’t quite right because there is a connection between η and its derivative; they are not absolutely independent, because $\eta(t)$ must be zero at both t_1 and t_2 .

“The method of solving all problems in the calculus of variations always uses the same general principle. You make the shift in the thing you want to vary (as we did by adding η); you look at the first-order terms; *then* you always arrange things in such a form that you get an integral of the form ‘some kind of stuff times the shift (η),’ but with no other derivatives (no $d\eta/dt$). It must be rearranged so it is always ‘something’ times η . You will see the great value of that in a minute. (There are formulas that tell you how to do this in some cases without actually calculating, but they are not general enough to be worth bothering about; the best way is to calculate it out this way.)

“How can I rearrange the term in $d\eta/dt$ to make it have an η ? I can do that by integrating by parts. It turns out that the whole trick of the calculus of variations consists of writing down the variation of S and then integrating by parts so that the derivatives of η disappear. It is always the same in every problem in which derivatives appear.

“You remember the general principle for integrating by parts. If you have any function f times $d\eta/dt$ integrated with respect to t , you write down the derivative of ηf :

$$\frac{d}{dt} (\eta f) = \eta \frac{df}{dt} + f \frac{d\eta}{dt}.$$

The integral you want is over the last term, so

$$\int f \frac{d\eta}{dt} dt = \eta f - \int \eta \frac{df}{dt} dt.$$

“In our formula for δS , the function f is m times dx/dt ; therefore, I have the following formula for δS .

$$\delta S = m \frac{dx}{dt} \eta(t) \Big|_{t_1}^{t_2} - \int_{t_1}^{t_2} \frac{d}{dt} \left(m \frac{dx}{dt} \right) \eta(t) dt - \int_{t_1}^{t_2} V'(\underline{x}) \eta(t) dt.$$

The first term must be evaluated at the two limits t_1 and t_2 . Then I must have the integral from the rest of the integration by parts. The last term is brought down without change.

“Now comes something which always happens—the integrated part disappears. (In fact, if the integrated part does not disappear, you restate the principle, adding conditions to make sure it does!) We have already said that η must be zero at both ends of the path, because the principle is that the action is a minimum provided that the varied curve begins and ends at the chosen points. The condition is that

$\eta(t_1) = 0$, and $\eta(t_2) = 0$. So the integrated term is zero. We collect the other terms together and obtain this:

$$\delta S = \int_{t_1}^{t_2} \left[-m \frac{d^2 x}{dt^2} - V'(x) \right] \eta(t) dt.$$

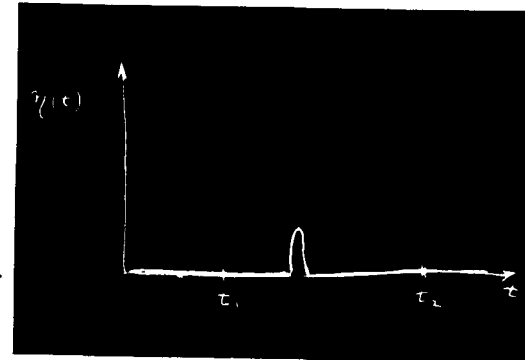
The variation in S is now the way we wanted it—there is the stuff in brackets, say F , all multiplied by $\eta(t)$ and integrated from t_1 to t_2 .

“We have that an integral of something or other times $\eta(t)$ is always zero:

$$\int F(t) \eta(t) dt = 0.$$

I have some function of t ; I multiply it by $\eta(t)$; and I integrate it from one end to the other. And no matter what the η is, I get zero. That means that the function $F(t)$ is zero. That’s obvious, but anyway I’ll show you one kind of proof.

“Suppose that for $\eta(t)$ I took something which was zero for all t except right near one particular value. It stays zero until it gets to this t ,



then it blips up for a moment and blips right back down. When we do the integral of this η times any function F , the only place that you get anything other than zero was where $\eta(t)$ was blipping, and then you get the value of F at that place times the integral over the blip. The integral over the blip alone isn’t zero, but when multiplied by F it has to be; so the function F has to be zero where the blip was. But the blip was anywhere I wanted to put it, so F must be zero everywhere.

“We see that if our integral is zero for any η , then the coefficient of η must be zero. The action integral will be a minimum for the path that satisfies this complicated differential equation:

$$\left[-m \frac{d^2 x}{dt^2} - V'(x) \right] = 0.$$

It’s not really so complicated; you have seen it before. It is just $F = ma$. The first term is the mass times acceleration, and the second is the derivative of the potential energy, which is the force.

“So, for a conservative system at least, we have demonstrated that the principle of least action gives the right answer; it says that the path that has the minimum action is the one satisfying Newton’s law.

“One remark: I did not prove it was a *minimum*—maybe it’s a maximum. In fact, it doesn’t really have to be a minimum. It is quite analogous to what we found for the ‘principle of least time’ which we discussed in optics. There also, we said at first it was ‘least’ time. It turned out, however, that there were situations in which it wasn’t the *least* time. The fundamental principle was that for any *first-order variation* away from the optical path, the *change* in time was zero; it is the same story. What we really mean by ‘least’ is that the first-order change in the value of S , when you change the path, is zero. It is not necessarily a ‘minimum.’

“Next, I remark on some generalizations. In the first place, the thing can be done in three dimensions. Instead of just x , I would have x , y , and z as functions of t ; the action is more complicated. For three-dimensional motion, you have to use the complete kinetic energy— $(m/2)$ times the whole velocity squared. That is,

$$\text{KE} = \frac{m}{2} \left[\left(\frac{dx}{dt} \right)^2 + \left(\frac{dy}{dt} \right)^2 + \left(\frac{dz}{dt} \right)^2 \right].$$

Also, the potential energy is a function of x , y , and z . And what about the path? The path is some general curve in space, which is not so easily drawn, but the idea is the same. And what about the η ? Well, η can have three components. You could shift the paths in x , or in y , or in z —or you could shift in all three directions simultaneously. So η would be a vector. This doesn’t really complicate things too much, though. Since only the *first-order* variation has to be zero, we can do the calculation by three successive shifts. We can shift η only in the x -direction and

say that coefficient must be zero. We get one equation. Then we shift it in the y -direction and get another. And in the z -direction and get another. Or, of course, in any order that you want. Anyway, you get three equations. And, of course, Newton's law is really three equations in the three dimensions—one for each component. I think that you can practically see that it is bound to work, but we will leave you to show for yourself that it will work for three dimensions. Incidentally, you could use any coordinate system you want, polar or otherwise, and get Newton's laws appropriate to that system right off by seeing what happens if you have the shift η in radius, or in angle, etc.

"Similarly, the method can be generalized to any number of particles. If you have, say, two particles with a force between them, so that there is a mutual potential energy, then you just add the kinetic energy of both particles and take the potential energy of the mutual interaction. And what do you vary? You vary the paths of *both* particles. Then, for two particles moving in three dimensions, there are six equations. You can vary the position of particle 1 in the x -direction, in the y -direction, and in the z -direction, and similarly for particle 2; so there are six equations. And that's as it should be. There are the three equations that determine the acceleration of particle 1 in terms of the force on it and three for the acceleration of particle 2, from the force on it. You follow the same game through, and you get Newton's law in three dimensions for any number of particles.

"I have been saying that we get Newton's law. That is not quite true, because Newton's law includes nonconservative forces like friction. Newton said that ma is equal to any F . But the principle of least action only works for *conservative* systems—where all forces can be gotten from a potential function. You know, however, that on a microscopic level—on the deepest level of physics—there are no nonconservative forces. Nonconservative forces, like friction, appear only because we neglect microscopic complications—there are just too many particles to analyze. But the *fundamental* laws *can* be put in the form of a principle of least action.

"Let me generalize still further. Suppose we ask what happens if the particle moves relativistically. We did not get the right relativistic equation of motion; $F = ma$ is only right nonrelativistically. The question is: Is there a corresponding principle of least action for the relativistic case? There is. The formula in the case of relativity is the following:

$$S = -m_0 c^2 \int_{t_1}^{t_2} \sqrt{1 - v^2/c^2} dt - q \int_{t_1}^{t_2} [\phi(x, y, z, t) - \mathbf{v} \cdot \mathbf{A}(x, y, z, t)] dt.$$

The first part of the action integral is the rest mass m_0 times c^2 times the integral of a function of velocity, $\sqrt{1 - v^2/c^2}$. Then instead of just the potential energy, we have an integral over the scalar potential ϕ and over \mathbf{v} times the vector potential \mathbf{A} . Of course, we are then including only electromagnetic forces. All electric and magnetic fields are given in terms of ϕ and \mathbf{A} . This action function gives the complete theory of relativistic motion of a single particle in an electromagnetic field.

"Of course, wherever I have written \mathbf{v} , you understand that before you try to figure anything out, you must substitute dx/dt for v_x and so on for the other components. Also, you put the point along the path at time t , $x(t), y(t), z(t)$ where I wrote simply x, y, z . Properly, it is only after you have made those replacements for the \mathbf{v} 's that you have the formula for the action for a relativistic particle. I will leave to the more ingenious of you the problem to demonstrate that this action formula does, in fact, give the correct equations of motion for relativity. May I suggest you do it first without the \mathbf{A} , that is, for no magnetic field? Then you should get the components of the equation of motion, $d\mathbf{p}/dt = -q \nabla \phi$, where, you remember, $\mathbf{p} = m\mathbf{v}/\sqrt{1 - v^2/c^2}$.

"It is much more difficult to include also the case with a vector potential. The variations get much more complicated. But in the end, the force term does come out equal to $q(\mathbf{E} + \mathbf{v} \times \mathbf{B})$, as it should. But I will leave that for you to play with.

"I would like to emphasize that in the general case, for instance in the relativistic formula, the action integrand no longer has the form of the kinetic energy

minus the potential energy. That's only true in the nonrelativistic approximation. For example, the term $m_0 c^2 \sqrt{1 - v^2/c^2}$ is not what we have called the kinetic energy. The question of what the action should be for any particular case must be determined by some kind of trial and error. It is just the same problem as determining what are the laws of motion in the first place. You just have to fiddle around with the equations that you know and see if you can get them into the form of the principle of least action.

"One other point on terminology. The function that is integrated over time to get the action S is called the *Lagrangian*, \mathcal{L} , which is a function only of the velocities and positions of particles. So the principle of least action is also written

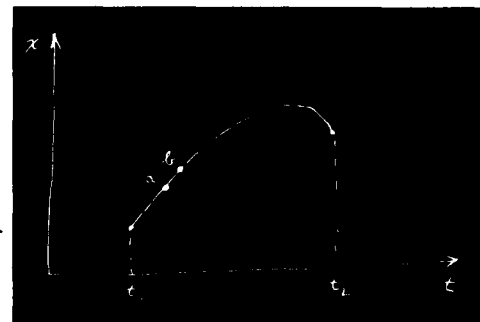
$$S = \int_{t_1}^{t_2} \mathcal{L}(x_i, v_i) dt,$$

where by x_i and v_i are meant all the components of the positions and velocities. So if you hear someone talking about the 'Lagrangian,' you know they are talking about the function that is used to find S . For relativistic motion in an electromagnetic field

$$\mathcal{L} = -m_0 c^2 \sqrt{1 - v^2/c^2} - q(\phi + \mathbf{v} \cdot \mathbf{A}).$$

"Also, I should say that S is not really called the 'action' by the most precise and pedantic people. It is called 'Hamilton's first principal function.' Now I hate to give a lecture on 'the-principle-of-least-Hamilton's-first-principal-function.' So I call it 'the action.' Also, more and more people are calling it the action. You see, historically something else which is not quite as useful was called the action, but I think it's more sensible to change to a newer definition. So now you too will call the new function the action, and pretty soon everybody will call it by that simple name.

"Now I want to say some things on this subject which are similar to the discussions I gave about the principle of least time. There is quite a difference in the characteristic of a law which says a certain integral from one place to another is a minimum—which tells something about the whole path—and of a law which says that as you go along, there is a force that makes it accelerate. The second way tells how you inch your way along the path, and the other is a grand statement about the whole path. In the case of light, we talked about the connection of these two. Now, I would like to explain why it is true that there are differential laws when there is a least action principle of this kind. The reason is the following: Consider the actual path in space and time. As before, let's take only one dimension, so we can plot the graph of x as a function of t . Along the true path, S is a minimum. Let's suppose that we have the true path and that it goes through some point a in space and time, and also through another nearby point b .



Now if the entire integral from t_1 to t_2 is a minimum, it is also necessary that the integral along the little section from a to b is also a minimum. It can't be that the part from a to b is a little bit more. Otherwise you could just fiddle with just that piece of the path and make the whole integral a little lower.

"So every subsection of the path must also be a minimum. And this is true no matter how short the subsection. Therefore, the principle that the whole path gives a minimum can be stated also by saying that an infinitesimal section of path also has a curve such that it has a minimum action. Now if we take a short enough section of path—between two points a and b very close together—how the potential varies from one place to another far away is not the important thing, because you are staying almost in the same place over the whole little piece of the path. The only thing that you have to discuss is the first-order change in the potential. The answer can only depend on the derivative of the potential and not on the potential everywhere. So the statement about the gross property of the whole path becomes a statement of what happens for a short section of the path—a differential statement. And this differential statement only involves the derivatives of the potential, that is, the force at a point. That's the qualitative explanation of the relation between the gross law and the differential law.

"In the case of light we also discussed the question: How does the particle find the right path? From the differential point of view, it is easy to understand. Every moment it gets an acceleration and knows only what to do at that instant. But all your instincts on cause and effect go haywire when you say that the particle decides to take the path that is going to give the minimum action. Does it 'smell' the neighboring paths to find out whether or not they have more action? In the case of light, when we put blocks in the way so that the photons could not test all the paths, we found that they couldn't figure out which way to go, and we had the phenomenon of diffraction.

"Is the same thing true in mechanics? Is it true that the particle doesn't just 'take the right path' but that it looks at all the other possible trajectories? And if by having things in the way, we don't let it look, that we will get an analog of diffraction? The miracle of it all is, of course, that it does just that. That's what the laws of quantum mechanics say. So our principle of least action is incompletely stated. It isn't that a particle takes the path of least action but that it smells all the paths in the neighborhood and chooses the one that has the least action by a method analogous to the one by which light chose the shortest time. You remember that the way light chose the shortest time was this: If it went on a path that took a different amount of time, it would arrive at a different phase. And the total amplitude at some point is the sum of contributions of amplitude for all the different ways the light can arrive. All the paths that give wildly different phases don't add up to anything. But if you can find a whole sequence of paths which have phases almost all the same, then the little contributions will add up and you get a reasonable total amplitude to arrive. The important path becomes the one for which there are many nearby paths which give the same phase.

"It is just exactly the same thing for quantum mechanics. The complete quantum mechanics (for the nonrelativistic case and neglecting electron spin) works as follows: The probability that a particle starting at point 1 at the time t_1 will arrive at point 2 at the time t_2 is the square of a probability amplitude. The total amplitude can be written as the sum of the amplitudes for each possible path—for each way of arrival. For every $x(t)$ that we could have—for every possible imaginary trajectory—we have to calculate an amplitude. Then we add them all together. What do we take for the amplitude for each path? Our action integral tells us what the amplitude for a single path ought to be. The amplitude is proportional to some constant times $e^{iS/\hbar}$, where S is the action for that path. That is, if we represent the phase of the amplitude by a complex number, the phase angle is S/\hbar . The action S has dimensions of energy times time, and Planck's constant \hbar has the same dimensions. It is the constant that determines when quantum mechanics is important.

"Here is how it works: Suppose that for all paths, S is very large compared to \hbar . One path contributes a certain amplitude. For a nearby path, the phase is quite different, because with an enormous S even a small change in S means a completely different phase—because \hbar is so tiny. So nearby paths will normally cancel their effects out in taking the sum—except for one region, and that is when a path and a nearby path all give the same phase in the first approximation (more precisely, the same action within \hbar). Only those paths will be the important ones. So in the limiting case in which Planck's constant \hbar goes to zero, the correct quantum-mechanical laws can be summarized by simply saying: 'Forget about all these probability amplitudes. The particle does go on a special path, namely, that one for which S does not vary in the first approximation.' That's the relation between the principle of least action and quantum mechanics. The fact that quantum mechanics can be formulated in this way was discovered in 1942 by a student of that same teacher, Bader, I spoke of at the beginning of this lecture. [Quantum mechanics was originally formulated by giving a differential equation for the amplitude (Schrödinger) and also by some other matrix mathematics (Heisenberg).]

"Now I want to talk about other minimum principles in physics. There are many very interesting ones. I will not try to list them all now but will only describe one more. Later on, when we come to a physical phenomenon which has a nice minimum principle, I will tell about it then. I want now to show that we can de-

scribe electrostatics, not by giving a differential equation for the field, but by saying that a certain integral is a maximum or a minimum. First, let's take the case where the charge density is known everywhere, and the problem is to find the potential ϕ everywhere in space. You know that the answer should be

$$\nabla^2 \phi = -\rho/\epsilon_0.$$

But another way of stating the same thing is this: Calculate the integral U^* , where

$$U^* = \frac{\epsilon_0}{2} \int (\nabla \phi)^2 dV - \int \rho \phi dV,$$

which is a volume integral to be taken over all space. This thing is a minimum for the correct potential distribution $\phi(x, y, z)$.

"We can show that the two statements about electrostatics are equivalent. Let's suppose that we pick any function ϕ . We want to show that when we take for ϕ the correct potential $\underline{\phi}$, plus a small deviation f , then in the first order, the change in U^* is zero. So we write

$$\phi = \underline{\phi} + f.$$

The $\underline{\phi}$ is what we are looking for, but we are making a variation of it to find what it has to be so that the variation of U^* is zero to first order. For the first part of U^* , we need

$$(\nabla \phi)^2 = (\nabla \underline{\phi})^2 + 2 \nabla \underline{\phi} \cdot \nabla f + (\nabla f)^2.$$

The only first-order term that will vary is

$$2 \nabla \underline{\phi} \cdot \nabla f.$$

In the second term of the quantity U^* , the integrand is

$$\rho \phi = \rho \underline{\phi} + \rho f,$$

whose variable part is ρf . So, keeping only the variable parts, we need the integral

$$\Delta U^* = \int (\epsilon_0 \nabla \underline{\phi} \cdot \nabla f - \rho f) dV.$$

"Now, following the old general rule, we have to get the darn thing all clear of derivatives of f . Let's look at what the derivatives are. The dot product is

$$\frac{\partial \underline{\phi}}{\partial x} \frac{\partial f}{\partial x} + \frac{\partial \underline{\phi}}{\partial y} \frac{\partial f}{\partial y} + \frac{\partial \underline{\phi}}{\partial z} \frac{\partial f}{\partial z},$$

which we have to integrate with respect to x , to y , and to z . Now here is the trick: to get rid of $\partial f / \partial x$ we integrate by parts with respect to x . That will carry the derivative over onto the $\underline{\phi}$. It's the same general idea we used to get rid of derivatives with respect to t . We use the equality

$$\int \frac{\partial \underline{\phi}}{\partial x} \frac{\partial f}{\partial x} dx = f \frac{\partial \underline{\phi}}{\partial x} - \int f \frac{\partial^2 \underline{\phi}}{\partial x^2} dx.$$

The integrated term is zero, since we have to make f zero at infinity. (That corresponds to making η zero at t_1 and t_2 . So our principle should be more accurately stated: U^* is less for the true ϕ than for any other $\phi(x, y, z)$ having the same values at infinity.) Then we do the same thing for y and z . So our integral ΔU^* is

$$\Delta U^* = \int (-\epsilon_0 \nabla^2 \underline{\phi} - \rho) f dV.$$

In order for this variation to be zero for any f , no matter what, the coefficient of f must be zero and, therefore,

$$\nabla^2 \underline{\phi} = -\rho/\epsilon_0.$$

We get back our old equation. So our 'minimum' proposition is correct.

"We can generalize our proposition if we do our algebra in a little different way. Let's go back and do our integration by parts without taking components. We start by looking at the following equality:

$$\nabla \cdot (f \nabla \underline{\phi}) = \nabla f \cdot \nabla \underline{\phi} + f \nabla^2 \underline{\phi}.$$

If I differentiate out the left-hand side, I can show that it is just equal to the right-hand side. Now we can use this equation to integrate by parts. In our integral ΔU^* , we replace $-\nabla \underline{\phi} \cdot \nabla f$ by $f \nabla^2 \underline{\phi} - \nabla \cdot (f \nabla \underline{\phi})$, which gets integrated over volume. The divergence term integrated over volume can be replaced by a surface integral:

$$\int \nabla \cdot (f \nabla \underline{\phi}) dV = \int f \nabla \underline{\phi} \cdot \underline{n} da.$$

Since we are integrating over all space, the surface over which we are integrating is at infinity. There, f is zero and we get the same answer as before.

"Only now we see how to solve a problem when we *don't* know where all the charges are. Suppose that we have conductors with charges spread out on them in some way. We can still use our minimum principle if the potentials of all the conductors are fixed. We carry out the integral for U^* only in the space outside of all conductors. Then, since we can't vary $\underline{\phi}$ on the conductor, f is zero on all those surfaces, and the surface integral

$$\int f \nabla \underline{\phi} \cdot \underline{n} da$$

is still zero. The remaining volume integral

$$\Delta U^* = \int (-\epsilon_0 \nabla^2 \underline{\phi} - \rho \underline{\phi}) f dV$$

is only to be carried out in the spaces between conductors. Of course, we get Poisson's equation again,

$$\nabla^2 \underline{\phi} = -\rho/\epsilon_0.$$

So we have shown that our original integral U^* is also a minimum if we evaluate it over the space outside of conductors all at fixed potentials (that is, such that any trial $\phi(x, y, z)$ must equal the given potential of the conductors when x, y, z is a point on the surface of a conductor).

"There is an interesting case when the only charges are on conductors. Then

$$U^* = \frac{\epsilon_0}{2} \int (\nabla \phi)^2 dV.$$

Our minimum principle says that in the case where there are conductors set at certain given potentials, the potential between them adjusts itself so that integral U^* is least. What is this integral? The term $\nabla \phi$ is the electric field, so the integral is the electrostatic energy. The true field is the one, of all those coming from the gradient of a potential, with the minimum total energy.

"I would like to use this result to calculate something particular to show you that these things are really quite practical. Suppose I take two conductors in the form of a cylindrical condenser.

The inside conductor has the potential V , and the outside is at the potential zero. Let the radius of the inside conductor be a and that of the outside, b . Now we can suppose *any* distribution of potential between the two. If we use the *correct* ϕ , and calculate $\epsilon_0/2 \int (\nabla \phi)^2 dV$, it should be the energy of the system, $\frac{1}{2} CV^2$.



So we can also calculate C by our principle. But if we use a wrong distribution of potential and try to calculate the capacity C by this method, we will get a capacity that is too big, since V is specified. Any assumed potential ϕ that is not the exactly correct one will give a fake C that is larger than the correct value. But if my false ϕ is any rough approximation, the C will be a good approximation, because the error in C is second order in the error in ϕ .

“Suppose I don’t know the capacity of a cylindrical condenser. I can use this principle to find it. I just guess at the potential function ϕ until I get the lowest C . Suppose, for instance, I pick a potential that corresponds to a constant field. (You know, of course, that the field isn’t really constant here; it varies as $1/r$.) A field which is constant means a potential which goes linearly with distance. To fit the conditions at the two conductors, it must be

$$\phi = V \left(1 - \frac{r-a}{b-a} \right).$$

This function is V at $r = a$, zero at $r = b$, and in between has a constant slope equal to $-V/(b-a)$. So what one does to find the integral U^* is multiply the square of this gradient by $\epsilon_0/2$ and integrate over all volume. Let’s do this calculation for a cylinder of unit length. A volume element at the radius r is $2\pi r dr$. Doing the integral, I find that my first try at the capacity gives

$$\frac{1}{2} C V^2 (\text{first try}) = \frac{\epsilon_0}{2} \int_a^b \frac{V^2}{(b-a)^2} 2\pi r dr.$$

The integral is easy; it is just

$$\pi V^2 \left(\frac{b+a}{b-a} \right).$$

So I have a formula for the capacity which is not the true one but is an approximate job:

$$\frac{C}{2\pi\epsilon_0} = \frac{b+a}{2(b-a)}.$$

It is, naturally, different from the correct answer $C = 2\pi\epsilon_0/\ln(b/a)$, but it’s not too bad. Let’s compare it with the right answer for several values of b/a . I have computed out the answers in this table:

$\frac{b}{a}$	$\frac{C_{\text{true}}}{2\pi\epsilon_0}$	$\frac{C \text{ (first approx.)}}{2\pi\epsilon_0}$
2	1.4423	1.500
4	0.721	0.833
10	0.434	0.612
100	0.267	0.51
1.5	2.4662	2.50
1.1	10.492070	10.500000

Even when b/a is as big as 2—which gives a pretty big variation in the field compared with a linearly varying field—I get a pretty fair approximation. The answer is, of course, a little too high, as expected. The thing gets much worse if you have a tiny wire inside a big cylinder. Then the field has enormous variations and if you represent it by a constant, you’re not doing very well. With $b/a = 100$, we’re off by nearly a factor of two. Things are much better for small b/a . To take the opposite extreme, when the conductors are not very far apart—say $b/a = 1.1$ —then the constant field is a pretty good approximation, and we get the correct value for C to within a tenth of a percent.

“Now I would like to tell you how to improve such a calculation. (Of course, you *know* the right answer for the cylinder, but the method is the same for some other odd shapes, where you may not know the right answer.) The next step is to try a better approximation to the unknown true ϕ . For example, we might try a

constant plus an exponential ϕ , etc. But how do you know when you have a better approximation unless you know the true ϕ ? Answer: You calculate C ; the lowest C is the value nearest the truth. Let us try this idea out. Suppose that the potential is not linear but say quadratic in r —that the electric field is not constant but linear. The most *general* quadratic form that fits $\phi = 0$ at $r = b$ and $\phi = V$ at $r = a$ is

$$\phi = V \left[1 + \alpha \left(\frac{r-a}{b-a} \right) - (1 + \alpha) \left(\frac{r-a}{b-a} \right)^2 \right],$$

where α is any constant number. This formula is a little more complicated. It involves a quadratic term in the potential as well as a linear term. It is very easy to get the field out of it. The field is just

$$E = -\frac{d\phi}{dr} = -\frac{\alpha V}{b-a} + 2(1 + \alpha) \frac{(r-a)V}{(b-a)^2}.$$

Now we have to square this and integrate over volume. But wait a moment. What should I take for α ? I can take a parabola for the ϕ ; but what parabola? Here's what I do: Calculate the capacity with *an arbitrary* α . What I get is

$$\frac{C}{2\pi\epsilon_0} = \frac{a}{b-a} \left[\frac{b}{a} \left(\frac{\alpha^2}{6} + \frac{2\alpha}{3} + 1 \right) + \frac{1}{6} \alpha^2 + \frac{1}{3} \right].$$

It looks a little complicated, but it comes out of integrating the square of the field. Now I can pick my α . I know that the truth lies lower than anything that I am going to calculate, so whatever I put in for α is going to give me an answer too big. But if I keep playing with α and get the lowest possible value I can, that lowest value is nearer to the truth than any other value. So what I do next is to pick the α that gives the minimum value for C . Working it out by ordinary calculus, I get that the minimum C occurs for $\alpha = -2b/(b+a)$. Substituting that value into the formula, I obtain for the minimum capacity

$$\frac{C}{2\pi\epsilon_0} = \frac{b^2 + 4ab + a^2}{3(b^2 - a^2)}.$$

"I've worked out what this formula gives for C for various values of b/a . I call these numbers $C(\text{quadratic})$. Here is a table that compares $C(\text{quadratic})$ with the true C .

$\frac{b}{a}$	$\frac{C_{\text{true}}}{2\pi\epsilon_0}$	$\frac{C(\text{quadratic})}{2\pi\epsilon_0}$
2	1.4423	1.444
4	0.721	0.733
10	0.434	0.475
100	0.267	0.346
1.5	2.4662	2.4667
1.1	10.492070	10.492065

"For example, when the ratio of the radii is 2 to 1, I have 1.444, which is a very good approximation to the true answer, 1.4423. Even for larger b/a , it stays pretty good—it is much, much better than the first approximation. It is even fairly good—only off by 10 percent—when b/a is 10 to 1. But when it gets to be 100 to 1—well, things begin to go wild. I get that C is 0.346 instead of 0.267. On the other hand, for a ratio of radii of 1.5, the answer is excellent; and for a b/a of 1.1, the answer comes out 10.492065 instead of 10.492070. Where the answer should be good, it is very, very good.

"I have given these examples, first, to show the theoretical value of the principles of minimum action and minimum principles in general and, second, to show their practical utility—not just to calculate a capacity when we already know the answer. For any other shape, you can guess an approximate field with some unknown parameters like α and adjust them to get a minimum. You will get excellent numerical results for otherwise intractable problems."

A note added after the lecture

"I should like to add something that I didn't have time for in the lecture. (I always seem to prepare more than I have time to tell about.) As I mentioned earlier, I got interested in a problem while working on this lecture. I want to tell you what that problem is. Among the minimum principles that I could mention, I noticed that most of them sprang in one way or another from the least action principle of mechanics and electrodynamics. But there is also a class that does not. As an example, if currents are made to go through a piece of material obeying Ohm's law, the currents distribute themselves inside the piece so that the rate at which heat is generated is as little as possible. Also we can say (if things are kept isothermal) that the rate at which energy is generated is a minimum. Now, this principle also holds, according to classical theory, in determining even the distribution of velocities of the electrons inside a metal which is carrying a current. The distribution of velocities is not exactly the equilibrium distribution [Chapter 40, Vol. I; Eq. (40.6)] because they are drifting sideways. The new distribution can be found from the principle that it is the distribution for a given current for which the entropy developed per second by collisions is as small as possible. The true description of the electrons' behavior ought to be by quantum mechanics, however. The question is: Does the same principle of minimum entropy generation also hold when the situation is described quantum-mechanically? I haven't found out yet.

"The question is interesting academically, of course. Such principles are fascinating, and it is always worth while to try to see how general they are. But also from a more practical point of view, I *want* to know. I, with some colleagues, have published a paper in which we calculated by quantum mechanics approximately the electrical resistance felt by an electron moving through an ionic crystal like NaCl. [Feynman, Hellworth, Iddings, and Platzman, "Mobility of Slow Electrons in a Polar Crystal," *Phys. Rev.* **127**, 1004 (1962).] But if a minimum principle existed, we could use it to make the results much more accurate, just as the minimum principle for the capacity of a condenser permitted us to get such accuracy for that capacity even though we had only a rough knowledge of the electric field."

Solutions of Maxwell's Equations in Free Space

20-1 Waves in free space; plane waves

In Chapter 18 we had reached the point where we had the Maxwell equations in complete form. All there is to know about the classical theory of the electric and magnetic fields can be found in the four equations:

$$\begin{aligned} \text{I. } \nabla \cdot \mathbf{E} &= \frac{\rho}{\epsilon_0} & \text{II. } \nabla \times \mathbf{E} &= -\frac{\partial \mathbf{B}}{\partial t} \\ \text{III. } \nabla \cdot \mathbf{B} &= 0 & \text{IV. } c^2 \nabla \times \mathbf{B} &= \frac{\mathbf{j}}{\epsilon_0} + \frac{\partial \mathbf{E}}{\partial t} \end{aligned} \quad (20.1)$$

When we put all these equations together, a remarkable new phenomenon occurs: fields generated by moving charges can leave the sources and travel alone through space. We considered a special example in which an infinite current sheet is suddenly turned on. After the current has been on for the time t , there are uniform electric and magnetic fields extending out the distance ct from the source. Suppose that the current sheet lies in the yz -plane with a surface current density J going toward positive y . The electric field will have only a y -component, and the magnetic field, only a z -component. The magnitude of the field components is given by

$$E_y = cB_z = -\frac{J}{2\epsilon_0 c}, \quad (20.2)$$

for positive values of x less than ct . For larger x the fields are zero. There are, of course, similar fields extending the same distance from the current sheet in the negative x -direction. In Fig. 20-1 we show a graph of the magnitude of the fields as a function of x at the instant t . As time goes on, the "wavefront" at ct moves outward in x at the constant velocity c .

Now consider the following sequence of events. We turn on a current of unit strength for a while, then suddenly increase the current strength to three units, and hold it constant at this value. What do the fields look like then? We can see what the fields will look like in the following way. First, we imagine a current of unit strength that is turned on at $t = 0$ and left constant forever. The fields for positive x are then given by the graph in part (a) of Fig. 20-2. Next, we ask what would happen if we turn on a steady current of two units at the time t_1 .

The fields in this case will be twice as high as before, but will extend out in x only the distance $c(t - t_1)$, as shown in part (b) of the figure. When we add these two solutions, using the principle of superposition, we find that the sum of the two sources is a current of one unit for the time from zero to t_1 and a current of three units for times greater than t_1 . At the time t the fields will vary with x as shown in part (c) of Fig. 20-2.

Now let's take a more complicated problem. Consider a current which is turned on to one unit for a while, then turned up to three units, and later turned off to zero. What are the fields for such a current? We can find the solution in the same way—by adding the solutions of three separate problems. First, we find the fields for a step current of unit strength. (We have solved that problem already.) Next, we find the fields produced by a step current of two units. Finally, we solve for the fields of a step current of *minus* three units. When we add the three solutions, we will have a current which is one unit strong from $t = 0$ to some later time, say t_1 , then three units strong until a still later time t_2 , and then turned off—that

20-1 Waves in free space; plane waves

20-2 Three-dimensional waves

20-3 Scientific imagination

20-4 Spherical waves

References: Chapter 47, Vol. I: *Sound: The Wave Equation*
Chapter 28, Vol. I: *Electromagnetic Radiation*

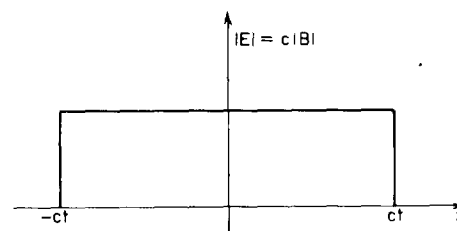


Fig. 20-1. The electric and magnetic field as a function of x at the time t after the current sheet is turned on.

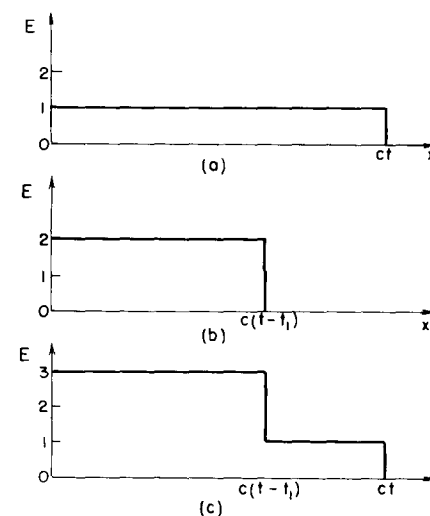


Fig. 20-2. The electric field of a current sheet. (a) One unit of current turned on at $t = 0$; (b) Two units of current turned on at $t = t_1$; (c) Superposition of (a) and (b).

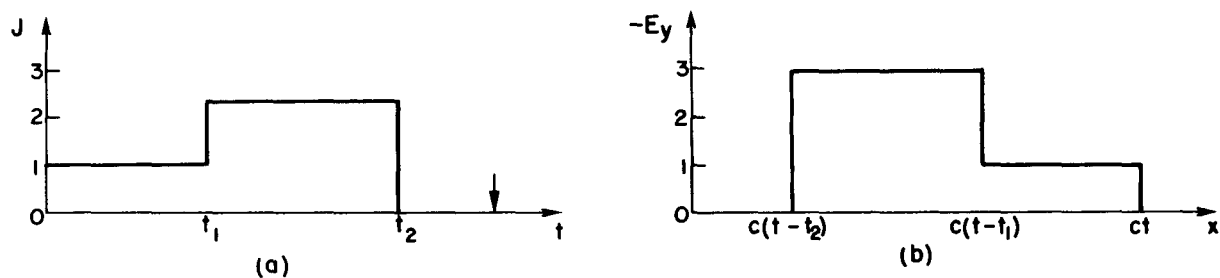


Fig. 20-3. If the current source strength varies as shown in (a), then at the time t shown by the arrow the electric field as a function of x is as shown in (b).

is, to zero. A graph of the current as a function of time is shown in Fig. 20-3(a). When we add the three solutions for the electric field, we find that its variation with x , at a given instant t , is as shown in Fig. 20-3(b). The field is an exact representation of the current. The field distribution in space is a nice graph of the current variation with time—only drawn backwards. As time goes on the whole picture moves outward at the speed c , so there is a little blob of field, travelling toward positive x , which contains a completely detailed memory of the history of all the current variations. If we were to stand miles away, we could tell from the variation of the electric or magnetic field exactly how the current had varied at the source.

You will also notice that long after all activity at the source has completely stopped and all charges and currents are zero, the block of field continues to travel through space. We have a distribution of electric and magnetic fields that exist independently of any charges or currents. That is the new effect that comes from the complete set of Maxwell's equations. If we want, we can give a complete mathematical representation of the analysis we have just done by writing that the electric field at a given place and a given time is proportional to the current at the source, only not at the *same* time, but at the *earlier* time $t - x/c$. We can write

$$E_y(t) = - \frac{J(t - x/c)}{2\epsilon_0 c}. \quad (20.3)$$

We have, believe it or not, already derived this same equation from another point of view in Vol. I, when we were dealing with the theory of the index of refraction. Then, we had to figure out what fields were produced by a thin layer of oscillating dipoles in a sheet of dielectric material with the dipoles set in motion by the electric field of an incoming electromagnetic wave. Our problem was to calculate the combined fields of the original wave and the waves radiated by the oscillating dipoles. How could we have calculated the fields generated by moving charges when we didn't have Maxwell's equations? At that time we took as our starting point (without any derivation) a formula for the radiation fields produced at large distances from an accelerating point charge. If you will look in Chapter 31 of Vol. I, you will see that Eq. (31.10) there is just the same as the Eq. (20.3) that we have just written down. Although our earlier derivation was correct only at large distances from the source, we see now that the same result continues to be correct even right up to the source.

We want now to look in a general way at the behavior of electric and magnetic fields in empty space far away from the sources, i.e., from the currents and charges. Very near the sources—near enough so that during the delay in transmission, the source has not had time to change much—the fields are very much the same as we have found in what we called the electrostatic or magnetostatic cases. If we go out to distances large enough so that the delays become important, however, the nature of the fields can be radically different from the solutions we have found. In a sense, the fields begin to take on a character of their own when they have gone a long way from all the sources. So we can begin by discussing the behavior of the fields in a region where there are no currents or charges.

Suppose we ask: What kind of fields can there be in regions where ρ and \mathbf{j} are both zero? In Chapter 18 we saw that the physics of Maxwell's equations could also be expressed in terms of differential equations for the scalar and vector potentials:

$$\nabla^2 \phi - \frac{1}{c^2} \frac{\partial^2 \phi}{\partial t^2} = - \frac{\rho}{\epsilon_0}, \quad (20.4)$$

$$\nabla^2 \mathbf{A} - \frac{1}{c^2} \frac{\partial^2 \mathbf{A}}{\partial t^2} = - \frac{\mathbf{j}}{\epsilon_0 c^2}. \quad (20.5)$$

If ρ and \mathbf{j} are zero, these equations take on the simpler form

$$\nabla^2 \phi - \frac{1}{c^2} \frac{\partial^2 \phi}{\partial t^2} = 0, \quad (20.6)$$

$$\nabla^2 \mathbf{A} - \frac{1}{c^2} \frac{\partial^2 \mathbf{A}}{\partial t^2} = 0. \quad (20.7)$$

Thus in free space the scalar potential ϕ and each component of the vector potential \mathbf{A} all satisfy the same mathematical equation. Suppose we let ψ (psi) stand for any one of the four quantities ϕ, A_x, A_y, A_z ; then we want to investigate the general solutions of the following equation:

$$\nabla^2 \psi - \frac{1}{c^2} \frac{\partial^2 \psi}{\partial t^2} = 0. \quad (20.8)$$

This equation is called the three-dimensional wave equation—three-dimensional, because the function ψ may depend in general on x, y , and z , and we need to worry about variations in all three coordinates. This is made clear if we write out explicitly the three terms of the Laplacian operator:

$$\frac{\partial^2 \psi}{\partial x^2} + \frac{\partial^2 \psi}{\partial y^2} + \frac{\partial^2 \psi}{\partial z^2} - \frac{1}{c^2} \frac{\partial^2 \psi}{\partial t^2} = 0. \quad (20.9)$$

In free space, the electric fields \mathbf{E} and \mathbf{B} also satisfy the wave equation. For example, since $\mathbf{B} = \nabla \times \mathbf{A}$, we can get a differential equation for \mathbf{B} by taking the curl of Eq. (20.7). Since the Laplacian is a scalar operator, the order of the Laplacian and curl operations can be interchanged:

$$\nabla \times (\nabla^2 \mathbf{A}) = \nabla^2 (\nabla \times \mathbf{A}) = \nabla^2 \mathbf{B}.$$

Similarly, the order of the operations curl and $\partial/\partial t$ can be interchanged:

$$\nabla \times \frac{1}{c^2} \frac{\partial^2 \mathbf{A}}{\partial t^2} = \frac{1}{c^2} \frac{\partial^2}{\partial t^2} (\nabla \times \mathbf{A}) = \frac{1}{c^2} \frac{\partial^2 \mathbf{B}}{\partial t^2}.$$

Using these results, we get the following differential equation for \mathbf{B} :

$$\nabla^2 \mathbf{B} - \frac{1}{c^2} \frac{\partial^2 \mathbf{B}}{\partial t^2} = 0. \quad (20.10)$$

So each component of the magnetic field \mathbf{B} satisfies the three-dimensional wave equation. Similarly, using the fact that $\mathbf{E} = -\nabla \phi - d\mathbf{A}/dt$, it follows that the electric field \mathbf{E} in free space also satisfies the three-dimensional wave equation:

$$\nabla^2 \mathbf{E} - \frac{1}{c^2} \frac{\partial^2 \mathbf{E}}{\partial t^2} = 0. \quad (20.11)$$

All of our electromagnetic fields satisfy the same wave equation, Eq. (20.8). We might well ask: What is the most general solution to this equation? However, rather than tackling that difficult question right away, we will look first at what can be said in general about those solutions in which nothing varies in y and z . (Always do an easy case first so that you can see what is going to happen, and then you can go to the more complicated cases.) Let's suppose that the magnitudes

of the fields depend only upon x —that there are no *variations* of the fields with y and z . We are, of course, considering plane waves again. We should expect to get results something like those in the previous section. In fact, we will find precisely the same answers. You may ask: “Why do it all over again?” It is important to do it again, first, because we did not show that the waves we found were the most general solutions for plane waves, and second, because we found the fields only from a very particular kind of current source. We would like to ask now: What is the most general kind of one-dimensional wave there can be in free space? We cannot find that by seeing what happens for this or that particular source, but must work with greater generality. Also we are going to work this time with differential equations instead of with integral forms. Although we will get the same results, it is a way of practicing back and forth to show that it doesn’t make any difference which way you go. You should know how to do things every which way, because when you get a hard problem, you will often find that only one of the various ways is tractable.

We could consider directly the solution of the wave equation for some electromagnetic quantity. Instead, we want to start right from the beginning with Maxwell’s equations in free space so that you can see their close relationship to the electromagnetic waves. So we start with the equations in (20.1), setting the charges and currents equal to zero. They become

$$\begin{aligned} \text{I. } \nabla \cdot \mathbf{E} &= 0 \\ \text{II. } \nabla \times \mathbf{E} &= - \frac{\partial \mathbf{B}}{\partial t} \\ \text{III. } \nabla \cdot \mathbf{B} &= 0 \\ \text{IV. } c^2 \nabla \times \mathbf{B} &= \frac{\partial \mathbf{E}}{\partial t} \end{aligned} \tag{20.12}$$

We write the first equation out in components:

$$\nabla \cdot \mathbf{E} = \frac{\partial E_x}{\partial x} + \frac{\partial E_y}{\partial y} + \frac{\partial E_z}{\partial z} = 0. \tag{20.13}$$

We are assuming that there are no variations with y and z , so the last two terms are zero. This equation then tells us that

$$\frac{\partial E_x}{\partial x} = 0. \tag{20.14}$$

Its solution is that E_x , the component of the electric field in the x -direction, is a constant in space. If you look at IV in (20.12), supposing no \mathbf{B} -variation in y and z either, you can see that E_x is also constant in time. Such a field could be the steady DC field from some charged condenser plates a long distance away. We are not interested now in such an uninteresting static field; we are at the moment interested only in dynamically varying fields. For *dynamic* fields, $E_x = 0$.

We have then the important result that for the propagation of plane waves in any direction, *the electric field must be at right angles to the direction of propagation*. It can, of course, still vary in a complicated way with the coordinate x .

The transverse \mathbf{E} -field can always be resolved into two components, say the y -component and the z -component. So let’s first work out a case in which the electric field has only one transverse component. We’ll take first an electric field that is always in the y -direction, with zero z -component. Evidently, if we solve this problem we can also solve for the case where the electric field is always in the z -direction. The general solution can always be expressed as the superposition of two such fields.

How easy our equations now get. The only component of the electric field that is not zero is E_y , and all derivatives—except those with respect to x —are zero. The rest of Maxwell’s equations then become quite simple.

Let's look next at the second of Maxwell's equations [II of Eq. (20.12)]. Writing out the components of the curl \mathbf{E} , we have

$$(\nabla \times \mathbf{E})_x = \frac{\partial E_z}{\partial y} - \frac{\partial E_y}{\partial z} = 0,$$

$$(\nabla \times \mathbf{E})_y = \frac{\partial E_x}{\partial z} - \frac{\partial E_z}{\partial x} = 0,$$

$$(\nabla \times \mathbf{E})_z = \frac{\partial E_y}{\partial x} - \frac{\partial E_x}{\partial y} = \frac{\partial E_y}{\partial x}.$$

The x -component of $\nabla \times \mathbf{E}$ is zero because the derivatives with respect to y and z are zero. The y -component is also zero; the first term is zero because the derivative with respect to z is zero, and the second term is zero because E_z is zero. The only components of the curl of \mathbf{E} that is not zero is the z -component, which is equal to $\partial E_y / \partial x$. Setting the three components of $\nabla \times \mathbf{E}$ equal to the corresponding components of $-\partial \mathbf{B} / \partial t$, we can conclude the following:

$$\frac{\partial B_x}{\partial t} = 0, \quad \frac{\partial B_y}{\partial t} = 0. \quad (20.15)$$

$$\frac{\partial B_z}{\partial t} = -\frac{\partial E_y}{\partial x}. \quad (20.16)$$

Since the x -component of the magnetic field and the y -component of the magnetic field both have zero time derivatives, these two components are just constant fields and correspond to the magnetostatic solutions we found earlier. Somebody may have left some permanent magnets near where the waves are propagating. We will ignore these constant fields and set B_x and B_y equal to zero.

Incidentally, we would already have concluded that the x -component of \mathbf{B} should be zero for a different reason. Since the divergence of \mathbf{B} is zero (from the third Maxwell equation), applying the same arguments we used above for the electric field, we would conclude that the longitudinal component of the magnetic field can have no variation with x . Since we are ignoring such uniform fields in our wave solutions, we would have set B_x equal to zero. In plane electromagnetic waves the \mathbf{B} -field, as well as the \mathbf{E} -field, must be directed at right angles to the direction of propagation.

Equation (20.16) gives us the additional proposition that if the electric field has only a y -component, the magnetic field will have only a z -component. So \mathbf{E} and \mathbf{B} are at right angles to each other. This is exactly what happened in the special wave we have already considered.

We are now ready to use the last of Maxwell's equations for free space [IV of Eq. (20.12)]. Writing out the components, we have

$$\begin{aligned} c^2(\nabla \times \mathbf{B})_x &= c^2 \frac{\partial B_z}{\partial y} - c^2 \frac{\partial B_y}{\partial z} = \frac{\partial E_x}{\partial t}, \\ c^2(\nabla \times \mathbf{B})_y &= c^2 \frac{\partial B_x}{\partial z} - c^2 \frac{\partial B_z}{\partial x} = \frac{\partial E_y}{\partial t}, \\ c^2(\nabla \times \mathbf{B})_z &= c^2 \frac{\partial B_y}{\partial x} - c^2 \frac{\partial B_x}{\partial y} = \frac{\partial E_z}{\partial t}. \end{aligned} \quad (20.17)$$

Of the six derivatives of the components of \mathbf{B} , only the term $\partial B_z / \partial x$ is not equal to zero. So the three equations give us simply

$$-c^2 \frac{\partial B_z}{\partial x} = \frac{\partial E_y}{\partial t}. \quad (20.18)$$

The result of all our work is that only one component each of the electric and magnetic fields is not zero, and that these components must satisfy Eqs. (20.16) and (20.18). The two equations can be combined into one if we differentiate the first with respect to x and the second with respect to t ; the left-hand sides of the

two equations will then be the same (except for the factor c^2). So we find that E_y satisfies the equation

$$\frac{\partial^2 E_y}{\partial x^2} - \frac{1}{c^2} \frac{\partial^2 E_y}{\partial t^2} = 0. \quad (20.19)$$

We have seen the same differential equation before, when we studied the propagation of sound. It is the wave equation for one-dimensional waves.

You should note that in the process of our derivation we have found something *more* than is contained in Eq. (20.11). Maxwell's equations have given us the further information that electromagnetic waves have field components only at right angles to the direction of the wave propagation.

Let's review what we know about the solutions of the one-dimensional wave equation. If any quantity ψ satisfies the one-dimensional wave equation

$$\frac{\partial^2 \psi}{\partial x^2} - \frac{1}{c^2} \frac{\partial^2 \psi}{\partial t^2} = 0, \quad (20.20)$$

then one possible solution is a function $\psi(x, t)$ of the form

$$\psi(x, t) = f(x - ct), \quad (20.21)$$

that is, some function of the *single* variable $(x - ct)$. The function $f(x - ct)$ represents a "rigid" pattern in x which travels toward positive x at the speed c (see Fig. 20-4). For example, if the function f has a maximum when its argument is zero, then for $t = 0$ the maximum of ψ will occur at $x = 0$. At some later time, say $t = 10$, ψ will have its maximum at $x = 10c$. As time goes on, the maximum moves toward positive x at the speed c .

Sometimes it is more convenient to say that a solution of the one-dimensional wave equation is a function of $(t - x/c)$. However, this is saying the same thing, because any function of $(t - x/c)$ is also a function of $(x - ct)$:

$$F(t - x/c) = F\left[-\frac{x - ct}{c}\right] = f(x - ct).$$

Let's show that $f(x - ct)$ is indeed a solution of the wave equation. Since it is a function of only one variable—the variable $(x - ct)$ —we will let f' represent the derivative of f with respect to its variable and f'' represent the second derivative of f . Differentiating Eq. (20.21) with respect to x , we have

$$\frac{\partial \psi}{\partial x} = f'(x - ct),$$

since the derivative of $(x - ct)$ with respect to x is 1. The second derivative of ψ with respect to x is clearly

$$\frac{\partial^2 \psi}{\partial x^2} = f''(x - ct). \quad (20.22)$$

Taking derivatives of ψ with respect to t , we find

$$\begin{aligned} \frac{\partial \psi}{\partial t} &= f'(x - ct)(-c), \\ \frac{\partial^2 \psi}{\partial t^2} &= +c^2 f''(x - ct) \end{aligned} \quad (20.23)$$

We see that ψ does indeed satisfy the one-dimensional wave equation.

You may be wondering: "If I have the wave equation, how do I know that I should take $f(x - ct)$ as a solution? I don't like this backward method. Isn't there some *forward* way to find the solution?" Well, one good forward way is to know the solution. It is possible to "cook up" an apparently forward mathematical argument, especially because we know what the solution is supposed to be, but with an equation as simple as this we don't have to play games. Soon you will get so that when you see Eq. (20.20), you nearly simultaneously see

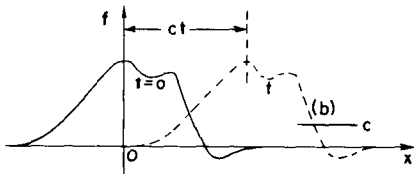


Fig. 20-4. The function $f(x - ct)$ represents a constant "shape" that travels toward positive x with the speed c .

$\psi = f(x - ct)$ as a solution. (Just as now when you see the integral of $x^2 dx$, you know right away that the answer is $x^3/3$.)

Actually you should also see a little more. Not only is any function of $(x - ct)$ a solution, but any function of $(x + ct)$ is also a solution. Since the wave equation contains only c^2 , changing the sign of c makes no difference. In fact, the *most general* solution of the one-dimensional wave equation is the sum of two arbitrary functions, one of $(x - ct)$ and the other of $(x + ct)$:

$$\psi = f(x - ct) + g(x + ct). \quad (20.24)$$

The first term represents a wave travelling toward positive x , and the second term an arbitrary wave travelling toward negative x . The general solution is the superposition of two such waves both existing at the same time.

We will leave the following amusing question for you to think about. Take a function ψ of the following form:

$$\psi = \cos kx \cos kct.$$

This equation isn't in the form of a function of $(x - ct)$ or of $(x + ct)$. Yet you can easily show that this function is a solution of the wave equation by direct substitution into Eq. (20.20). How can we then say that the general solution is of the form of Eq. (20.24)?

Applying our conclusions about the solution of the wave equation to the y -component of the electric field, E_y , we conclude that E_y can vary with x in any arbitrary fashion. However, the fields which do exist can always be considered as the sum of two patterns. One wave is sailing through space in one direction with speed c , with an associated magnetic field perpendicular to the electric field; another wave is travelling in the opposite direction with the same speed. Such waves correspond to the electromagnetic waves that we know about—light, radio-waves, infrared radiation, ultraviolet radiation, x-rays, and so on. We have already discussed the radiation of light in great detail in Vol. I. Since everything we learned there applies to any electromagnetic wave, we don't need to consider in great detail here the behavior of these waves.

We should perhaps make a few further remarks on the question of the polarization of the electromagnetic waves. In our solution we chose to consider the special case in which the electric field has only a y -component. There is clearly another solution for waves travelling in the plus or minus x -direction, with an electric field which has only a z -component. Since Maxwell's equations are linear, the general solution for one-dimensional waves propagating in the x -direction is the sum of waves of E_y and waves of E_z . This general solution is summarized in the following equations:

$$\begin{aligned} \mathbf{E} &= (0, E_y, E_z) \\ E_y &= f(x - ct) + g(x + ct) \\ E_z &= F(x - ct) + G(x + ct) \\ \mathbf{B} &= (0, B_y, B_z) \\ cB_z &= f(x - ct) - g(x + ct) \\ cB_y &= -F(x - ct) + G(x + ct). \end{aligned} \quad (20.25)$$

Such electromagnetic waves have an \mathbf{E} -vector whose direction is not constant but which gyrates around in some arbitrary way in the yz -plane. At every point the magnetic field is always perpendicular to the electric field and to the direction of propagation.

If there are only waves travelling in one direction, say the positive x -direction, there is a simple rule which tells the relative orientation of the electric and magnetic fields. The rule is that the cross product $\mathbf{E} \times \mathbf{B}$ —which is, of course, a vector at right angles to both \mathbf{E} and \mathbf{B} —points in the direction in which the wave is travelling. If \mathbf{E} is rotated into \mathbf{B} by a right-hand screw, the screw points in the direction of the wave velocity. (We shall see later that the vector $\mathbf{E} \times \mathbf{B}$ has a special physical significance: it is a vector which describes the flow of energy in an electromagnetic field.)

20-2 Three-dimensional waves

We want now to turn to the subject of three-dimensional waves. We have already seen that the vector \mathbf{E} satisfies the wave equation. It is also easy to arrive at the same conclusion by arguing directly from Maxwell's equations. Suppose we start with the equation

$$\nabla \times \mathbf{E} = - \frac{\partial \mathbf{B}}{\partial t}$$

and take the curl of both sides:

$$\nabla \times (\nabla \times \mathbf{E}) = - \frac{\partial}{\partial t} (\nabla \times \mathbf{B}). \quad (20.26)$$

You will remember that the curl of the curl of any vector can be written as the sum of two terms, one involving the divergence and the other the Laplacian,

$$\nabla \times (\nabla \times \mathbf{E}) = \nabla(\nabla \cdot \mathbf{E}) - \nabla^2 \mathbf{E}.$$

In free space, however, the divergence of \mathbf{E} is zero, so only the Laplacian term remains. Also, from the fourth of Maxwell's equations in free space [Eq. (20.12)] the time derivative of $c^2 \nabla \times \mathbf{B}$ is the second derivative of \mathbf{E} with respect to t :

$$c^2 \frac{\partial}{\partial t} (\nabla \times \mathbf{B}) = \frac{\partial^2 \mathbf{E}}{\partial t^2}.$$

Equation (20.26) then becomes

$$\nabla^2 \mathbf{E} = \frac{1}{c^2} \frac{\partial^2 \mathbf{E}}{\partial t^2},$$

which is the three-dimensional wave equation. Written out in all its glory, this equation is, of course,

$$\frac{\partial^2 \mathbf{E}}{\partial x^2} + \frac{\partial^2 \mathbf{E}}{\partial y^2} + \frac{\partial^2 \mathbf{E}}{\partial z^2} - \frac{1}{c^2} \frac{\partial^2 \mathbf{E}}{\partial t^2} = 0. \quad (20.27)$$

How shall we find the general wave solution? The answer is that all the solutions of the three-dimensional wave equation can be represented as a superposition of the one-dimensional solutions we have already found. We obtained the equation for waves which move in the x -direction by supposing that the field did not depend on y and z . Obviously, there are other solutions in which the fields do not depend on x and z , representing waves going in the y -direction. Then there are solutions which do not depend on x and y , representing waves travelling in the z -direction. Or in general, since we have written our equations in vector form, the three-dimensional wave equation can have solutions which are plane waves moving in any direction at all. Again, since the equations are linear, we may have simultaneously as many plane waves as we wish, travelling in as many different directions. Thus the most general solution of the three-dimensional wave equation is a superposition of all sorts of plane waves moving in all sorts of directions.

Try to imagine what the electric and magnetic fields look like at present in the space in this lecture room. First of all, there is a steady magnetic field; it comes from the currents in the interior of the earth—that is, the earth's steady magnetic field. Then there are some irregular, nearly static electric fields produced perhaps by electric charges generated by friction as various people move about in their

chairs and rub their coat sleeves against the chair arms. Then there are other magnetic fields produced by oscillating currents in the electrical wiring—fields which vary at a frequency of 60 cycles per second, in synchronism with the generator at Boulder Dam. But more interesting are the electric and magnetic fields varying at much higher frequencies. For instance, as light travels from window to floor and wall to wall, there are little wiggles of the electric and magnetic fields moving along at 186,000 miles per second. Then there are also infrared waves travelling from the warm foreheads to the cold blackboard. And we have forgotten the ultraviolet light, the x-rays, and the radiowaves travelling through the room.

Flying across the room are electromagnetic waves which carry music of a jazz band. There are waves modulated by a series of impulses representing pictures of events going on in other parts of the world, or of imaginary aspirins dissolving in imaginary stomachs. To demonstrate the reality of these waves it is only necessary to turn on electronic equipment that converts these waves into pictures and sounds.

If we go into further detail to analyze even the smallest wiggles, there are tiny electromagnetic waves that have come into the room from enormous distances. There are now tiny oscillations of the electric field, whose crests are separated by a distance of one foot, that have come from millions of miles away, transmitted to the earth from the Mariner II space craft which has just passed Venus. Its signals carry summaries of information it has picked up about the planets (information obtained from electromagnetic waves that travelled from the planet to the space craft).

There are very tiny wiggles of the electric and magnetic fields that are waves which originated billions of light years away—from galaxies in the remotest corners of the universe. That this is true has been found by “filling the room with wires”—by building antennas as large as this room. Such radiowaves have been detected from places in space beyond the range of the greatest optical telescopes. Even they, the optical telescopes, are simply gatherers of electromagnetic waves. What we call the stars are only inferences, inferences drawn from the only physical reality we have yet gotten from them—from a careful study of the unendingly complex undulations of the electric and magnetic fields reaching us on earth.

There is, of course, more: the fields produced by lightning miles away, the fields of the charged cosmic ray particles as they zip through the room, and more, and more. What a complicated thing is the electric field in the space around you! Yet it always satisfies the three-dimensional wave equation.

20-3 Scientific imagination

I have asked you to imagine these electric and magnetic fields. What do you do? Do you know how? How do *I* imagine the electric and magnetic field? What do *I* actually see? What are the demands of scientific imagination? Is it any different from trying to imagine that the room is full of invisible angels? No, it is not like imagining invisible angels. It requires a much higher degree of imagination to understand the electromagnetic field than to understand invisible angels. Why? Because to make invisible angels understandable, all I have to do is to alter their properties *a little bit*—I make them slightly visible, and then I can see the shapes of their wings, and bodies, and halos. Once I succeed in imagining a visible angel, the abstraction required—which is to take almost invisible angels and imagine them completely invisible—is relatively easy. So you say, “Professor, please give me an approximate description of the electromagnetic waves, even though it may be slightly inaccurate, so that I too can see them as well as I can see almost invisible angels. Then I will modify the picture to the necessary abstraction.”

I’m sorry I can’t do that for you. I don’t know how. I have no picture of this electromagnetic field that is in any sense accurate. I have known about the electromagnetic field a long time—I was in the same position 25 years ago that you are now, and I have had 25 years more of experience thinking about these wiggling waves. When I start describing the magnetic field moving through space, I speak of the *E*- and *B* fields and wave my arms and you may imagine that I can see them.

I'll tell you what I see. I see some kind of vague shadowy, wiggling lines—here and there is an E and B written on them somehow, and perhaps some of the lines have arrows on them—an arrow here or there which disappears when I look too closely at it. When I talk about the fields swishing through space, I have a terrible confusion between the symbols I use to describe the objects and the objects themselves. I cannot really make a picture that is even nearly like the true waves. So if you have some difficulty in making such a picture, you should not be worried that your difficulty is unusual.

Our science makes terrific demands on the imagination. The degree of imagination that is required is much more extreme than that required for some of the ancient ideas. The modern ideas are much harder to imagine. We use a lot of tools, though. We use mathematical equations and rules, and make a lot of pictures. What I realize now is that when I talk about the electromagnetic field in space, I see some kind of a superposition of all of the diagrams which I've ever seen drawn about them. I don't see little bundles of field lines running about because it worries me that if I ran at a different speed the bundles would disappear. I don't even always see the electric and magnetic fields because sometimes I think I should have made a picture with the vector potential and the scalar potential, for those were perhaps the more physically significant things that were wiggling.

Perhaps the only hope, you say, is to take a mathematical view. Now what is a mathematical view? From a mathematical view, there is an electric field vector and a magnetic field vector at every point in space; that is, there are six numbers associated with every point. Can you imagine six numbers associated with each point in space? That's too hard. Can you imagine even *one* number associated with every point? I cannot! I can imagine such a thing as the temperature at every point in space. That seems to be understandable. There is a hotness and coldness that varies from place to place. But I honestly do not understand the idea of a *number* at every point.

So perhaps we should put the question: Can we represent the electric field by something more like a temperature, say like the displacement of a piece of jello? Suppose that we were to begin by imagining that the world was filled with thin jello and that the fields represented some distortion—say a stretching or twisting—of the jello. Then we could visualize the field. After we “see” what it is like we could abstract the jello away. For many years that's what people tried to do. Maxwell, Ampere, Faraday, and others tried to understand electromagnetism this way. (Sometimes they called the abstract jello “ether.”) But it turned out that the attempt to imagine the electromagnetic field in that way was really standing in the way of progress. We are unfortunately limited to abstractions, to using instruments to detect the field, to using mathematical symbols to describe the field, etc. But nevertheless, in some sense the fields are real, because after we are all finished fiddling around with mathematical equations—with or without making pictures and drawings or trying to visualize the thing—we can still make the instruments detect the signals from Mariner II and find out about galaxies a billion miles away, and so on.

The whole question of imagination in science is often misunderstood by people in other disciplines. They try to test our imagination in the following way. They say, “Here is a picture of some people in a situation. What do you imagine will happen next?” When we say, “I can't imagine,” they may think we have a weak imagination. They overlook the fact that whatever we are *allowed* to imagine in science must be *consistent with everything else we know*: that the electric fields and the waves we talk about are not just some happy thoughts which we are free to make as we wish, but ideas which must be consistent with all the laws of physics we know. We can't allow ourselves to seriously imagine things which are obviously in contradiction to the known laws of nature. And so our kind of imagination is quite a difficult game. One has to have the imagination to think of something that has never been seen before, never been heard of before. At the same time the thoughts are restricted in a strait jacket, so to speak, limited by the conditions that come from our knowledge of the way nature really is. The problem of creating

something which is new, but which is consistent with everything which has been seen before, is one of extreme difficulty.

While I'm on this subject I want to talk about whether it will ever be possible to imagine *beauty* that we can't *see*. It is an interesting question. When we look at a rainbow, it looks beautiful to us. Everybody says, "Ooh, a rainbow." (You see how scientific I am. I am afraid to say something is beautiful unless I have an experimental way of defining it.) But how would we describe a rainbow if we were blind? We *are* blind when we measure the infrared reflection coefficient of sodium chloride, or when we talk about the frequency of the waves that are coming from some galaxy that we can't see—we make a diagram, we make a plot. For instance, for the rainbow, such a plot would be the intensity of radiation vs. wavelength measured with a spectrophotometer for each direction in the sky. Generally, such measurements would give a curve that was rather flat. Then some day, someone would discover that for certain conditions of the weather, and at certain angles in the sky, the spectrum of intensity as a function of wavelength would behave strangely; it would have a bump. As the angle of the instrument was varied only a little bit, the maximum of the bump would move from one wavelength to another. Then one day the physical review of the blind men might publish a technical article with the title "The Intensity of Radiation as a Function of Angle under Certain Conditions of the Weather." In this article there might appear a graph such as the one in Fig. 20-5. The author would perhaps remark that at the larger angles there was more radiation at long wavelengths, whereas for the smaller angles the maximum in the radiation came at shorter wavelengths. (From our point of view, we would say that the light at 40° is predominantly green and the light at 42° is predominantly red.)

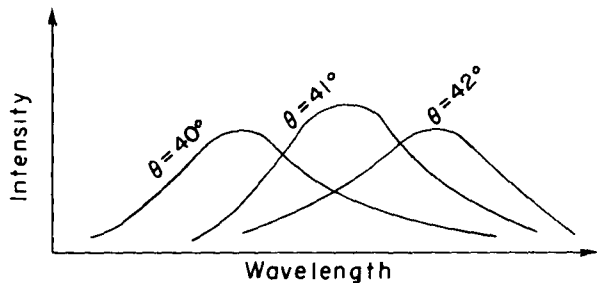


Fig. 20-5. The intensity of electromagnetic waves as a function of wavelength for three angles (measured from the direction opposite the sun), observed only with certain meteorological conditions.

Now do we find the graph of Fig. 20-5 beautiful? It contains much more detail than we apprehend when we look at a rainbow, because our eyes cannot see the exact details in the shape of a spectrum. The eye, however, finds the rainbow beautiful. Do we have enough imagination to see in the spectral curves the same beauty we see when we look directly at the rainbow? I don't know.

But suppose I have a graph of the reflection coefficient of a sodium chloride crystal as a function of wavelength in the infrared, and also as a function of angle. I would have a representation of how it would look to my eyes if they could see in the infrared—perhaps some glowing, shiny "green," mixed with reflections from the surface in a "metallic red." That would be a beautiful thing, but I don't know whether I can ever look at a graph of the reflection coefficient of NaCl measured with some instrument and say that it has the same beauty.

On the other hand, even if we cannot see beauty in particular measured results, we *can* already claim to see a certain beauty in the equations which describe general physical laws. For example, in the wave equation (20.9), there's something nice about the regularity of the appearance of the x , the y , the z , and the t . And this nice symmetry in appearance of the x , y , z , and t suggests to the mind still a greater beauty which has to do with the four dimensions, the possibility that space has four-dimensional symmetry, the possibility of analyzing that and the developments of the special theory of relativity. So there is plenty of intellectual beauty associated with the equations.

20-4 Spherical waves

We have seen that there are solutions of the wave equation which correspond to plane waves, and that any electromagnetic wave can be described as a superposition of many plane waves. In certain special cases, however, it is more convenient to describe the wave field in a different mathematical form. We would like to discuss now the theory of spherical waves—waves which correspond to spherical surfaces that are spreading out from some center. When you drop a stone into a lake, the ripples spread out in circular waves on the surface—they are two-dimensional waves. A spherical wave is a similar thing except that it spreads out in three dimensions.

Before we start describing spherical waves, we need a little mathematics. Suppose we have a function that depends only on the radial distance r from a certain origin—in other words, a function that is spherically symmetric. Let's call the function $\psi(r)$, where by r we mean

$$r = \sqrt{x^2 + y^2 + z^2},$$

the radial distance from the origin. In order to find out what functions $\psi(r)$ satisfy the wave equation, we will need an expression for the Laplacian of ψ . So we want to find the sum of the second derivatives of ψ with respect to x , y , and z . We will use the notation that $\psi'(r)$ represents the derivative of ψ with respect to r and $\psi''(r)$ represents the second derivative of ψ with respect to r .

First, we find the derivatives with respect to x . The first derivative is

$$\frac{\partial \psi(r)}{\partial x} = \psi'(r) \frac{\partial r}{\partial x}.$$

The second derivative of ψ with respect to x is

$$\frac{\partial^2 \psi}{\partial x^2} = \psi'' \left(\frac{\partial r}{\partial x} \right)^2 + \psi' \frac{\partial^2 r}{\partial x^2}.$$

We can evaluate the partial derivatives of r with respect to x from

$$\frac{\partial r}{\partial x} = \frac{x}{r}, \quad \frac{\partial^2 r}{\partial x^2} = \frac{1}{r} \left(1 - \frac{x^2}{r^2} \right).$$

So the second derivative of ψ with respect to x is

$$\frac{\partial^2 \psi}{\partial x^2} = \frac{x^2}{r^2} \psi'' + \frac{1}{r} \left(1 - \frac{x^2}{r^2} \right) \psi'. \quad (20.28)$$

Likewise,

$$\frac{\partial^2 \psi}{\partial y^2} = \frac{y^2}{r^2} \psi'' + \frac{1}{r} \left(1 - \frac{y^2}{r^2} \right) \psi', \quad (20.29)$$

$$\frac{\partial^2 \psi}{\partial z^2} = \frac{z^2}{r^2} \psi'' + \frac{1}{r} \left(1 - \frac{z^2}{r^2} \right) \psi'. \quad (20.30)$$

The Laplacian is the sum of these three derivatives. Remembering that $x^2 + y^2 + z^2 = r^2$, we get

$$\nabla^2 \psi(r) = \psi''(r) + \frac{2}{r} \psi'(r). \quad (20.31)$$

It is often more convenient to write this equation in the following form:

$$\nabla^2 \psi = \frac{1}{r} \frac{d^2}{dr^2} (r\psi). \quad (20.32)$$

If you carry out the differentiation indicated in Eq. (20.32), you will see that the right-hand side is the same as in Eq. (20.31).

If we wish to consider spherically symmetric fields which can propagate as spherical waves, our field quantity must be a function of both r and t . Suppose

we ask, then, what functions $\psi(r, t)$ are solutions of the three-dimensional wave equation

$$\nabla^2 \psi(r, t) - \frac{1}{c^2} \frac{\partial^2}{\partial t^2} \psi(r, t) = 0. \quad (20.33)$$

Since $\psi(r, t)$ depends only on the spatial coordinates through r , we can use the equation for the Laplacian we found above, Eq. (20.32). To be precise, however, since ψ is also a function of t , we should write the derivatives with respect to r as partial derivatives. Then the wave equation becomes

$$\frac{1}{r} \frac{\partial^2}{\partial r^2} (r\psi) - \frac{1}{c^2} \frac{\partial^2}{\partial t^2} \psi = 0.$$

We must now solve this equation, which appears to be much more complicated than the plane wave case. But notice that if we multiply this equation by r , we get

$$\frac{\partial^2}{\partial r^2} (r\psi) - \frac{1}{c^2} \frac{\partial^2}{\partial t^2} (r\psi) = 0. \quad (20.34)$$

This equation tells us that the function $r\psi$ satisfies the one-dimensional wave equation in the variable r . Using the general principle which we have emphasized so often, that the same equations always have the same solutions, we know that if $r\psi$ is a function only of $(r - ct)$ then it will be a solution of Eq. (20.34). So we know that spherical waves must have the form

$$r\psi(r, t) = f(r - ct).$$

Or, as we have seen before, we can equally well say that $r\psi$ can have the form

$$r\psi = f(t - r/c).$$

Dividing by r , we find that the field quantity ψ (whatever it may be) has the following form:

$$\psi = \frac{f(t - r/c)}{r}. \quad (20.35)$$

Such a function represents a general spherical wave travelling outward from the origin at the speed c . If we forget about the r in the denominator for a moment, the amplitude of the wave as a function of the distance from the origin at a given time has a certain shape that travels outward at the speed c . The factor r in the denominator, however, says that the amplitude of the wave decreases in proportion to $1/r$ as the wave propagates. In other words, unlike a plane wave in which the amplitude remains constant as the wave runs along, in a spherical wave the amplitude steadily decreases, as shown in Fig. 20-6. This effect is easy to understand from a simple physical argument.

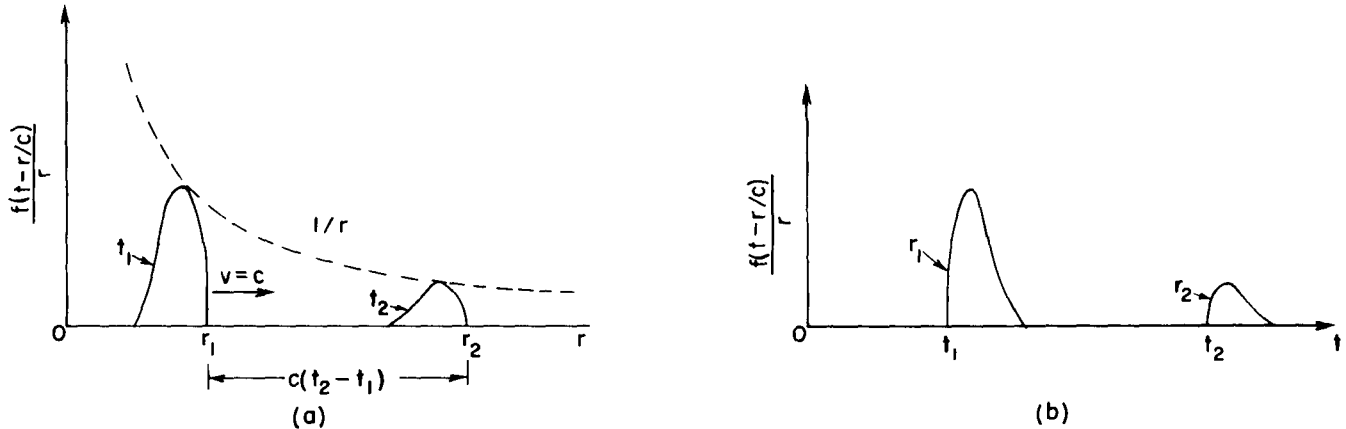


Fig. 20-6. A spherical wave $\psi = f(t - r/c)/r$. (a) ψ as a function of r for $t = t_1$ and the same wave for the later time t_2 . (b) ψ as a function of t for $r = r_1$ and the same wave seen at r_2 .

We know that the energy density in a wave depends on the square of the wave amplitude. As the wave spreads, its energy is spread over larger and larger areas proportional to the radial distance squared. If the total energy is conserved, the energy density must fall as $1/r^2$, and the amplitude of the wave must decrease as $1/r$. So Eq. (20.35) is the “reasonable” form for a spherical wave.

We have disregarded the second possible solution to the one-dimensional wave equation:

$$r\psi = g(t + r/c),$$

or

$$\psi = \frac{g(t + r/c)}{r}.$$

This also represents a spherical wave, but one which travels *inward* from large r toward the origin.

We are now going to make a special assumption. We say, without any demonstration whatever, that the waves generated by a source are only the waves which go *outward*. Since we know that waves are caused by the motion of charges, we want to think that the waves proceed outward from the charges. It would be rather strange to imagine that before charges were set in motion, a spherical wave started out from infinity and arrived at the charges just at the time they began to move. That is a possible solution, but experience shows that when charges are accelerated the waves travel outward from the charges. Although Maxwell’s equations would allow either possibility, we will put in an *additional fact*—based on experience—that only the outgoing wave solution makes “physical sense.”

We should remark, however, that there is an interesting consequence to this additional assumption: we are removing the symmetry with respect to time that exists in Maxwell’s equations. The original equations for \mathbf{E} and \mathbf{B} , and also the wave equations we derived from them, have the property that if we change the sign of t , the equation is unchanged. These equations say that for every solution corresponding to a wave going in one direction there is an equally valid solution for a wave travelling in the opposite direction. Our statement that we will consider only the outgoing spherical waves is an important additional assumption. (A formulation of electrodynamics in which this additional assumption is avoided has been carefully studied. Surprisingly, in many circumstances it does *not* lead to physically absurd conclusions, but it would take us too far astray to discuss these ideas just now. We will talk about them a little more in Chapter 28.)

We must mention another important point. In our solution for an outgoing wave, Eq. (20.35), the function ψ is infinite at the origin. That is somewhat peculiar. We would like to have a wave solution which is smooth everywhere. Our solution must represent physically a situation in which there is some source at the origin. In other words, we have inadvertently made a mistake. We have not solved the free wave equation (20.33) *everywhere*; we have solved Eq. (20.33) with zero on the right everywhere, except at the origin. Our mistake crept in because some of the steps in our derivation are not “legal” when $r = 0$.

Let’s show that it is easy to make the same kind of mistake in an electrostatic problem. Suppose we want a solution of the equation for an electrostatic potential in free space, $\nabla^2\phi = 0$. The Laplacian is equal to zero, because we are assuming that there are no charges anywhere. But what about a spherically symmetric solution to this equation—that is, some function ϕ that depends only on r . Using the formula of Eq. (20.32) for the Laplacian, we have

$$\frac{1}{r} \frac{d^2}{dr^2} (r\phi) = 0.$$

Multiplying this equation by r , we have an equation which is readily integrated:

$$\frac{d^2}{dr^2} (r\phi) = 0.$$

If we integrate once with respect to r , we find that the first derivative of $r\phi$ is a

constant, which we may call a :

$$\frac{d}{dr}(r\phi) = a.$$

Integrating again, we find that $r\phi$ is of the form

$$r\phi = ar + b,$$

where b is another constant of integration. So we have found that the following ϕ is a solution for the electrostatic potential in free space:

$$\phi = a + \frac{b}{r}.$$

Something is evidently wrong. In the region where there are no electric charges, we know the solution for the electrostatic potential: the potential is everywhere a constant. That corresponds to the first term in our solution. But we also have the second term, which says that there is a contribution to the potential that varies as one over the distance from the origin. We know, however, that such a potential corresponds to a point charge at the origin. So, although we thought we were solving for the potential in free space, our solution also gives the field for a point source at the origin. Do you see the similarity between what happened now and what happened when we solved for a spherically symmetric solution to the wave equation? If there were really no charges or currents at the origin, there would not be spherical outgoing waves. The spherical waves must, of course, be produced by sources at the origin. In the next chapter we will investigate the connection between the outgoing electromagnetic waves and the currents and voltages which produce them.

Solutions of Maxwell's Equations with Currents and Charges

21-1 Light and electromagnetic waves

We saw in the last chapter that among their solutions, Maxwell's equations have waves of electricity and magnetism. These waves correspond to the phenomena of radio, light, x-rays, and so on, depending on the wavelength. We have already studied light in great detail in Vol. I. In this chapter we want to tie together the two subjects—we want to show that Maxwell's equations can indeed form the base for our earlier treatment of the phenomena of light.

When we studied light, we began by writing down an equation for the electric field produced by a charge which moves in any arbitrary way. That equation was

$$\mathbf{E} = \frac{q}{4\pi\epsilon_0} \left[\frac{\mathbf{e}_{r'}}{r'^2} + \frac{r'}{c} \frac{d}{dt} \left(\frac{\mathbf{e}_{r'}}{r'^2} \right) + \frac{1}{c^2} \frac{d^2}{dt^2} \mathbf{e}_{r'} \right], \quad (21.1)$$

$$c\mathbf{B} = \mathbf{e}_{r'} \times \mathbf{E}.$$

[See Eq. (28.3), Vol. I.]

If a charge moves in an arbitrary way, the electric field we would find *now* at some point depends only on the position and motion of the charge not now, but at an *earlier* time—at an instant which is earlier by the time it would take light, going at the speed c , to travel the distance r' from the charge to the field point. In other words, if we want the electric field at point (1) at the time t , we must calculate the location (2') of the charge and its motion at the time $(t - r'/c)$, where r' is the distance to the point (1) from the position of the charge (2') at the time $(t - r'/c)$. The prime is to remind you that r' is the so-called “retarded distance” from the point (2') to the point (1), and not the actual distance between point (2), the position of the charge at the time t , and the field point (1) (see Fig. 21-1). Note that we are using a different convention now for the *direction* of the unit vector \mathbf{e}_r . In Chapters 28 and 36 of Vol. I it was convenient to take \mathbf{r} (and hence \mathbf{e}_r) pointing *toward* the source. Now we are following the definition we took for Coulomb's law, in which \mathbf{r} is directed *from* the charge, at (2), *toward* the field point at (1). The only difference, of course, is that our new \mathbf{r} (and \mathbf{e}_r) are the negatives of the old ones.

We have also seen that if the velocity v of a charge is always much less than c , and if we consider only points at large distances from the charge, so that only the last term of Eq. (21.1) is important, the fields can also be written as

$$\mathbf{E} = - \frac{q}{4\pi\epsilon_0 c^2 r'} \left[\begin{array}{l} \text{acceleration of the charge at } (t - r'/c) \\ \text{projected at right angles to } \mathbf{r}' \end{array} \right], \quad (21.1')$$

and

$$c\mathbf{B} = \mathbf{e}_{r'} \times \mathbf{E}.$$

Let's look at what the complete equation, Eq. (21.1), says in a little more detail. The vector $\mathbf{e}_{r'}$ is the unit vector to point (1) from the retarded position (2'). The first term, then, is what we would expect for the Coulomb field of the charge at its retarded position—we may call this “the retarded Coulomb field.” The electric field depends inversely on the square of the distance and is directed away from the retarded position of the charge (that is, in the direction of $\mathbf{e}_{r'}$).

But that is only the first term. The other terms tell us that the laws of electricity do *not* say that all the fields are the same as the static ones, but just retarded (which is what people sometimes like to say). To the “retarded Coulomb field” we must

21-1 Light and electromagnetic waves

21-2 Spherical waves from a point source

21-3 The general solution of Maxwell's equations

21-4 The fields of an oscillating dipole

21-5 The potentials of a moving charge; the general solution of Liénard and Wiechert

21-6 The potentials for a charge moving with constant velocity; the Lorentz formula

Review: Chapter 28, Vol. I, *Electromagnetic Radiation*
Chapter 31, Vol. I, *The Origin of the Refractive Index*
Chapter 36, Vol. I, *Relativistic Effects in Radiation*

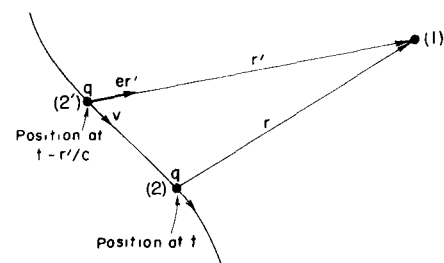


Fig. 21-1. The fields at (1) at the time t depend on the position (2') occupied by the charge q at the time $(t - r'/c)$.

add the other two terms. The second term says that there is a “correction” to the retarded Coulomb field which is the *rate of change* of the retarded Coulomb field multiplied by r'/c , the retardation delay. In a way of speaking, this term tends to *compensate* for the retardation in the first term. The first *two* terms correspond to computing the “retarded Coulomb field” and then extrapolating it toward the future by the amount r'/c , that is, *right up to the time t* ! The extrapolation is linear, as if we were to assume that the “retarded Coulomb field” would continue to change at the rate computed for the charge at the point (2'). If the field is changing slowly, the effect of the retardation is almost completely removed by the correction term, and the two terms together give us an electric field that is the “instantaneous Coulomb field”—that is, the Coulomb field of the charge at the point (2)—to a very good approximation.

Finally, there is a third term in Eq. (21.1) which is the second derivative of the unit vector $e_{r'}$. For our study of the phenomena of light, we made use of the fact that far away from the charge the first two terms went inversely as the square of the distance and, for large distances, became very weak in comparison to the last term, which decreases as $1/r$. So we concentrated entirely on the last term, and we showed that it is (again, for large distances) proportional to the component of the acceleration of the charge at right angles to the line of sight. (Also, for most of our work in Vol. I, we took the case in which the charges were moving nonrelativistically. We considered the relativistic effects in only one chapter, Chapter 36.)

Now we should try to connect the two things together. We have the Maxwell equations, and we have Eq. (21.1) for the field of a point charge. We should certainly ask whether they are equivalent. If we can deduce Eq. (21.1) from Maxwell's equations, we will really understand the connection between light and electromagnetism. To make this connection is the main purpose of this chapter.

It turns out that we won't quite make it—that the mathematical details get too complicated for us to carry through in all their gory details. But we will come close enough so that you should easily see how the connection could be made. The missing pieces will only be in the mathematical details. Some of you may find the mathematics in this chapter rather complicated, and you may not wish to follow the argument very closely. We think it is important, however, to make the connection between what you have learned earlier and what you are learning now, or at least to indicate how such a connection can be made. You will notice, if you look over the earlier chapters, that whenever we have taken a statement as a starting point for a discussion, we have carefully explained whether it is a new “assumption” that is a “basic law,” or whether it can ultimately be deduced from some other laws. We owe it to you in the spirit of these lectures to make the connection between light and Maxwell's equations. If it gets difficult in places, well, that's life—there is no other way.

21-2 Spherical waves from a point source

In Chapter 18 we found that Maxwell's equations could be solved by letting

$$\mathbf{E} = -\nabla\phi - \frac{\partial\mathbf{A}}{\partial t} \quad (21.2)$$

and

$$\mathbf{B} = \nabla \times \mathbf{A}, \quad (21.3)$$

where ϕ and \mathbf{A} must then be solutions of the equations

$$\nabla^2\phi - \frac{1}{c^2} \frac{\partial^2\phi}{\partial t^2} = -\frac{\rho}{\epsilon_0} \quad (21.4)$$

and

$$\nabla^2\mathbf{A} - \frac{1}{c^2} \frac{\partial^2\mathbf{A}}{\partial t^2} = -\frac{\mathbf{j}}{\epsilon_0 c^2}, \quad (21.5)$$

and must also satisfy the condition that

$$\nabla \cdot \mathbf{A} = -\frac{1}{c^2} \frac{\partial\phi}{\partial t}. \quad (21.6)$$

Now we will find the solution of Eqs. (21.4) and (21.5). To do that we have to find the solution ψ of the equation

$$\nabla^2 \psi - \frac{1}{c^2} \frac{\partial^2 \psi}{\partial t^2} = -s, \quad (21.7)$$

where s , which we call the source, is known. Of course, s corresponds to ρ/ϵ_0 and ψ to ϕ for Eq. (21.4), or s is $j_x/\epsilon_0 c^2$ if ψ is A_x , etc., but we want to solve Eq. (21.7) as a mathematical problem no matter what ψ and s are physically.

In places where ρ and j are zero—in what we have called “free” space—the potentials ϕ and A , and the fields E and B , all satisfy the three-dimensional wave equation without sources, whose mathematical form is

$$\nabla^2 \psi - \frac{1}{c^2} \frac{\partial^2 \psi}{\partial t^2} = 0. \quad (21.8)$$

In Chapter 20 we saw that solutions of this equation can represent waves of various kinds: plane waves in the x -direction, $\psi = f(t - x/c)$; plane waves in the y - or z -direction, or in any other direction; or spherical waves of the form

$$\psi(x, y, z, t) = \frac{f(t - r/c)}{r}. \quad (21.9)$$

(The solutions can be written in still other ways, for example cylindrical waves that spread out from an axis.)

We also remarked that, physically, Eq. (21.9) does not represent a wave in free space—that there must be charges at the origin to get the outgoing wave started. In other words, Eq. (21.9) is a solution of Eq. (21.8) everywhere except right near $r = 0$, where it must be a solution of the complete equation (21.7), including some sources. Let's see how that works. What kind of a source s in Eq. (21.7) would give rise to a wave like Eq. (21.9)?

Suppose we have the spherical wave of Eq. (21.9) and look at what is happening for very small r . Then the retardation $-r/c$ in $f(t - r/c)$ can be neglected—provided f is a smooth function—and ψ becomes

$$\psi = \frac{f(t)}{r} \quad (r \rightarrow 0). \quad (21.10)$$

So ψ is just like a Coulomb field for a charge at the origin that varies with time. That is, if we had a little lump of charge, limited to a very small region near the origin, with a density ρ , we know that

$$\phi = \frac{Q/4\pi\epsilon_0}{r},$$

where $Q = \int \rho dV$. Now we know that such a ϕ satisfies the equation

$$\nabla^2 \phi = -\frac{\rho}{\epsilon_0}.$$

Following the same mathematics, we would say that the ψ of Eq. (21.10) satisfies

$$\nabla^2 \psi = -s \quad (r \rightarrow 0), \quad (21.11)$$

where s is related to f by

$$f = \frac{S}{4\pi},$$

with

$$S = \int s dV.$$

The only difference is that in the general case, s , and therefore S , can be a function of time.

Now the important thing is that if ψ satisfies Eq. (21.11) for small r , it also satisfies Eq. (21.7). As we go very close to the origin, the $1/r$ dependence of ψ

causes the space derivatives to become very large. But the time derivatives keep their same values [They are just the time derivatives of $f(t)$.] So as r goes to zero, the term $\partial^2\psi/\partial t^2$ in Eq. (21.7) can be neglected in comparison with $\nabla^2\psi$, and Eq. (21.7) becomes equivalent to Eq. (21.11).

To summarize, then, if the source function $s(t)$ of Eq. (21.7) is localized at the origin and has the total strength

$$S(t) = \int s(t) dV, \quad (21.12)$$

the solution of Eq. (21.7) is

$$\psi(x, y, z, t) = \frac{1}{4\pi} \frac{S(t - r/c)}{r}. \quad (21.13)$$

The only effect of the term $\partial^2\psi/\partial t^2$ in Eq. (21.7) is to introduce the retardation $(t - r/c)$ in the Coulomb-like potential.

21-3 The general solution of Maxwell's equations

We have found the solution of Eq. (21.7) for a "point" source. The next question is: What is the solution for a spread-out source? That's easy; we can think of any source $s(x, y, z, t)$ as made up of the sum of many "point" sources, one for each volume element dV , and each with the source strength $s(x, y, z, t) dV$. Since Eq. (21.7) is linear, the resultant field is the superposition of the fields from all of such source elements.

Using the results of the preceding section [Eq. (21.13)] we know that the field $d\psi$ at the point (x_1, y_1, z_1) —or (1) for short—at the time t , from a source element $s dV$ at the point (x_2, y_2, z_2) —or (2) for short—is given by

$$d\psi(1, t) = \frac{s(2, t - r_{12}/c) dV_2}{4\pi r_{12}},$$

where r_{12} is the distance from (2) to (1). Adding the contributions from all the pieces of the source means, of course, doing an integral over all regions where $s \neq 0$; so we have

$$\psi(1, t) = \int \frac{s(2, t - r_{12}/c)}{4\pi r_{12}} dV_2. \quad (21.14)$$

That is, the field at (1) at the time t is the sum of all the spherical waves which leave the source elements at (2) at the times $(t - r_{12}/c)$. This is the solution of our wave equation for any set of sources.

We see now how to obtain a general solution for Maxwell's equations. If for ψ we mean the scalar potential ϕ , the source function s becomes ρ/ϵ_0 . Or we can let ψ represent any one of the three components of the vector potential A , replacing s by the corresponding component of $\mathbf{j}/\epsilon_0 c^2$. Thus, if we know the charge density $\rho(x, y, z, t)$ and the current density $\mathbf{j}(x, y, z, t)$ everywhere, we can immediately write down the solutions of Eqs. (21.4) and (21.5). They are

$$\phi(1, t) = \int \frac{\rho(2, t - r_{12}/c)}{4\pi\epsilon_0 r_{12}} dV_2 \quad (21.15)$$

and

$$A(1, t) = \int \frac{\mathbf{j}(2, t - r_{12}/c)}{4\pi\epsilon_0 c^2 r_{12}} dV_2. \quad (21.16)$$

The fields \mathbf{E} and \mathbf{B} can then be found by differentiating the potentials, using Eqs. (21.2) and (21.3). [Incidentally, it is possible to verify that the ϕ and A obtained from Eqs. (21.15) and (21.16) do satisfy the equality (21.6)]

We have solved Maxwell's equations. Given the currents and charges in any circumstance, we can find the potentials directly from these integrals and then differentiate and get the fields. So we have finished with the Maxwell theory. Also this permits us to close the ring back to our theory of light, because to connect with our earlier work on light, we need only calculate the electric field from a

moving charge. All that remains is to take a moving charge, calculate the potentials from these integrals, and then differentiate to find \mathbf{E} from $-\nabla\phi - \partial\mathbf{A}/\partial t$. We should get Eq. (21.1). It turns out to be lots of work, but that's the principle.

So here is the center of the universe of electromagnetism—the complete theory of electricity and magnetism, and of light; a complete description of the fields produced by any moving charges; and more. It is all here. Here is the structure built by Maxwell, complete in all its power and beauty. It is probably one of the greatest accomplishments of physics. To remind you of its importance, we will put it all together in a nice frame.

Maxwell's equations:

$$\nabla \cdot \mathbf{E} = \frac{\rho}{\epsilon_0}$$

$$\nabla \cdot \mathbf{B} = 0$$

$$\nabla \times \mathbf{E} = -\frac{\partial \mathbf{B}}{\partial t}$$

$$c^2 \nabla \times \mathbf{B} = \frac{\mathbf{j}}{\epsilon_0} + \frac{\partial \mathbf{E}}{\partial t}$$

Their solutions:

$$\mathbf{E} = -\nabla\phi - \frac{\partial \mathbf{A}}{\partial t}$$

$$\mathbf{B} = \nabla \times \mathbf{A}$$

$$\phi(1, t) = \int \frac{\rho(2, t - r_{12}/c)}{4\pi\epsilon_0 r_{12}} dV_2$$

$$\mathbf{A}(1, t) = \int \frac{\mathbf{j}(2, t - r_{12}/c)}{4\pi\epsilon_0 c^2 r_{12}} dV_2$$

21-4 The fields of an oscillating dipole

We have still not lived up to our promise to derive Eq. (21.1) for the electric field of a point charge in motion. Even with the results we already have, it is a relatively complicated thing to derive. We have not found Eq. (21.1) anywhere in the published literature except in Vol. I of these lectures.* So you can see that it is not easy to derive. (The fields of a moving charge have been written in many other forms that are equivalent, of course.) We will have to limit ourselves here just to showing that, in a few examples, Eqs. (21.15) and (21.16) give the same results as Eq. (21.1). First, we will show that Eq. (21.1) gives the correct fields with only the restriction that the motion of the charged particle is nonrelativistic. (Just this special case will take care of 90 percent, or more, of what we said about light.)

We consider a situation in which we have a blob of charge that is moving about in some way, in a small region, and we will find the fields far away. To put it another way, we are finding the field at any distance from a point charge that is shaking up and down in very small motion. Since light is usually emitted from neutral objects such as atoms, we will consider that our wiggling charge q is located near an equal and opposite charge at rest. If the separation between the centers of the charges is \mathbf{d} , the charges will have a dipole moment $\mathbf{p} = q\mathbf{d}$, which we take to be a function of time. Now we should expect that if we look at the fields close to the charges, we won't have to worry about the delay; the electric field will be exactly the same as the one we have calculated earlier for an electrostatic dipole

* The formula was worked out by R. P. Feynman, in about 1950, and given in some lectures as a good way of thinking about synchrotron radiation

—using, of course, the instantaneous dipole moment $\mathbf{p}(t)$. But if we go very far out, we ought to find a term in the field that goes as $1/r$ and depends on the acceleration of the charge perpendicular to the line of sight. Let's see if we get such a result.

We begin by calculating the vector potential \mathbf{A} , using Eq. (21.16). Suppose that our moving charge is in a small blob whose charge density is given by $\rho(x, y, z)$, and the whole thing is moving at any instant with the velocity \mathbf{v} . Then the current density $\mathbf{j}(x, y, z)$ will be equal to $\mathbf{v}\rho(x, y, z)$. It will be convenient to take our coordinate system so that the z -axis is in the direction of \mathbf{v} ; then the geometry of our problem is as shown in Fig. 21-2. We want the integral

$$\int \frac{\mathbf{j}(2, t - r_{12}/c)}{r_{12}} dV_2. \quad (21.17)$$

Now if the size of the charge-blob is really very small compared with r_{12} , we can set the r_{12} term in the denominator equal to r , the distance to the center of the blob, and take r outside the integral. Next, we are also going to set $r_{12} = r$ in the numerator, although that is not really quite right. It is not right because we should take \mathbf{j} at, say, the top of the blob at a slightly different time than we used for \mathbf{j} at the bottom of the blob. When we set $r_{12} = r$ in $\mathbf{j}(t - r_{12}/c)$, we are taking the current density for the whole blob at the same time $(t - r/c)$. That is an approximation that will be good only if the velocity v of the charge is much less than c . So we are making a nonrelativistic calculation. Replacing \mathbf{j} by $\rho\mathbf{v}$, the integral (21.17) becomes

$$\frac{1}{r} \int \mathbf{v}\rho(2, t - r/c) dV_2.$$

Since all the charge has the same velocity, this integral is just \mathbf{v}/r times the total charge q . But qv is just $\partial\mathbf{p}/\partial t$, the rate of change of the dipole moment—which is, of course, to be evaluated at the retarded time $(t - r/c)$. We will write it as $\dot{\mathbf{p}}(t - r/c)$. So we get for the vector potential

$$\mathbf{A}(1, t) = \frac{1}{4\pi\epsilon_0 c^2} \frac{\dot{\mathbf{p}}(t - r/c)}{r}. \quad (21.18)$$

Our result says that the current in a varying dipole produces a vector potential in the form of spherical waves whose source strength is $\dot{\mathbf{p}}/4\pi\epsilon_0 c^2$.

We can now get the magnetic field from $\mathbf{B} = \nabla \times \mathbf{A}$. Since \mathbf{p} is totally in the z -direction, \mathbf{A} has only a z -component; there are only two nonzero derivatives in the curl. So $B_x = \partial A_z / \partial y$ and $B_y = -\partial A_z / \partial x$. Let's first look at B_x :

$$B_x = \frac{\partial A_z}{\partial y} = \frac{1}{4\pi\epsilon_0 c^2} \frac{\partial}{\partial y} \frac{\dot{\mathbf{p}}(t - r/c)}{r}. \quad (21.19)$$

To carry out the differentiation, we must remember that $r = \sqrt{x^2 + y^2 + z^2}$, so

$$B_x = \frac{1}{4\pi\epsilon_0 c^2} \dot{\mathbf{p}}(t - r/c) \frac{\partial}{\partial y} \left(\frac{1}{r} \right) + \frac{1}{4\pi\epsilon_0 c^2} \frac{1}{r} \frac{\partial}{\partial y} \dot{\mathbf{p}}(t - r/c). \quad (21.20)$$

Remembering that $\partial r / \partial y = y/r$, the first term gives

$$- \frac{1}{4\pi\epsilon_0 c^2} \frac{y \dot{\mathbf{p}}(t - r/c)}{r^3}, \quad (21.21)$$

which drops off as $1/r^2$ like the fields of a static dipole (because y/r is constant for a given direction).

The second term in Eq. (21.20) gives us the new effects. Carrying out the differentiation, we get

$$- \frac{1}{4\pi\epsilon_0 c^2} \frac{y}{r^2} \ddot{\mathbf{p}}(t - r/c), \quad (21.22)$$

where $\ddot{\mathbf{p}}$ means, of course, the second derivative of \mathbf{p} with respect to t . This term,

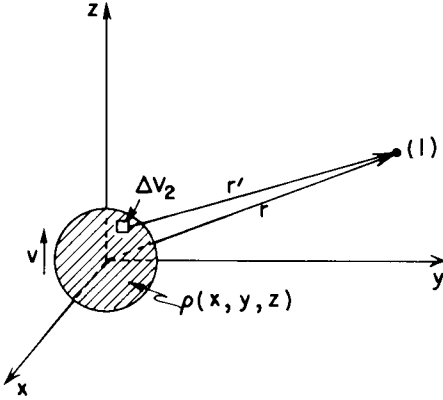


Fig. 21-2. The potentials at (1) are given by integrals over the charge density ρ .

which comes from differentiating the numerator, is responsible for radiation. First, it describes a field which decreases with distance only as $1/r$. Second, it depends on the *acceleration* of the charge. You can begin to see how we are going to get a result like Eq. (21.1'), which describes the radiation of light.

Let's examine in a little more detail how this radiation term comes about—it is such an interesting and important result. We start with the expression (21.18), which has a $1/r$ dependence and is therefore like a Coulomb potential, except for the delay term in the numerator. Why is it then that when we differentiate with respect to space coordinates to get the fields, we don't just get a $1/r^2$ field—with, of course, the corresponding time delays?

We can see why in the following way: Suppose that we let our dipole oscillate up and down in a sinusoidal motion. Then we would have

$$p = p_z = p_0 \sin \omega t$$

and

$$A_z = \frac{1}{4\pi\epsilon_0 c^2} \frac{\omega p_0 \cos \omega(t - r/c)}{r}.$$

If we plot a graph of A_z as a function of r at a given instant, we get the curve shown in Fig. 21-3. The peak amplitude decreases as $1/r$, but there is, in addition, an oscillation in space, bounded by the $1/r$ envelope. When we take the spatial derivatives, they will be proportional to the *slope* of the curve. From the figure we see that there are slopes much steeper than the slope of the $1/r$ curve itself. It is, in fact, evident that for a given frequency the peak slopes are proportional to the amplitude of the wave, which varies as $1/r$. So that explains the drop-off rate of the radiation term.

It all comes about because the variations *with time* at the source are translated into variations *in space* as the waves are propagated outward, and the magnetic fields depend on the *spatial* derivatives of the potential.

Let's go back and finish our calculation of the magnetic field. We have for B_x the two terms (21.21) and (21.22), so

$$B_x = \frac{1}{4\pi\epsilon_0 c^2} \left[-\frac{y\dot{p}(t - r/c)}{r^3} - \frac{y\ddot{p}(t - r/c)}{cr^2} \right].$$

With the same kind of mathematics, we get

$$B_y = \frac{1}{4\pi\epsilon_0 c^2} \left[\frac{x\dot{p}(t - r/c)}{r^3} + \frac{x\ddot{p}(t - r/c)}{cr^2} \right].$$

Or we can put it all together in a nice vector formula:

$$\mathbf{B} = \frac{1}{4\pi\epsilon_0 c^2} \frac{[\mathbf{p} + (r/c)\dot{\mathbf{p}}]_{t-r/c} \times \mathbf{r}}{r^3}. \quad (21.23)$$

Now let's look at this formula. First of all, if we go very far out in r , only the $\ddot{\mathbf{p}}$ term counts. The direction of \mathbf{B} is given by $\mathbf{p} \times \mathbf{r}$, which is at right angles to the radius \mathbf{r} and also at right angles to the acceleration, as in Fig. 21-4. Everything is coming out right; that is also the result we get from Eq. (21.1').

Now let's look at what we are not used to—at what happens closer in. In Section 14-9 we worked out the law of Biot and Savart for the magnetic field of an element of current. We found that a current element $\mathbf{j} dV$ contributes to the magnetic field the amount

$$d\mathbf{B} = \frac{1}{4\pi\epsilon_0 c^2} \frac{\mathbf{j} \times \mathbf{r}}{r^3} dV. \quad (21.24)$$

You see that this formula looks very much like the first term of Eq. (21.23), if we remember that \mathbf{p} is the current. But there is one difference. In Eq. (21.23), the current is to be evaluated at the time $(t - r/c)$, which doesn't appear in Eq. (21.24). Actually, however, Eq. (21.24) is still very good for small r , because the *second*

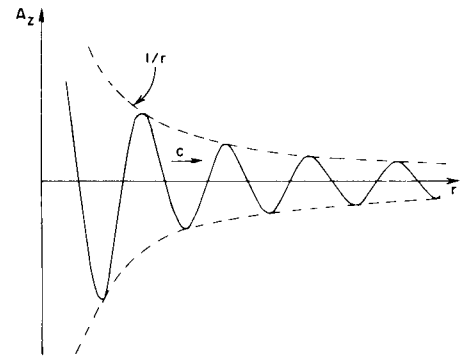


Fig. 21-3. The magnitude of \mathbf{A} as a function of r at the instant t for the spherical wave from an oscillating dipole.

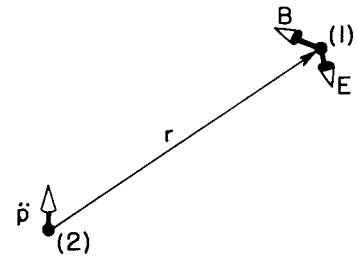


Fig. 21-4. The radiation fields \mathbf{B} and \mathbf{E} of an oscillating dipole.

term of Eq. (21.23) tends to cancel out the effect of the retardation in the first term. The two *together* give a result very near to Eq. (21.24) when r is small.

We can see that this way. When r is small, $(t - r/c)$ is not very different from t , so we can expand the bracket in Eq. (21.23) in a Taylor series. For the first term,

$$\mathbf{p}(t - r/c) = \mathbf{p}(t) - \frac{r}{c} \dot{\mathbf{p}}(t) + \text{etc.},$$

and to the same order in r/c ,

$$\dot{\mathbf{p}}(t - r/c) = \ddot{\mathbf{p}}(t).$$

When we take the sum, the two terms in $\dot{\mathbf{p}}$ cancel, and we are left with the *unretarded* current \mathbf{p} : that is, $\mathbf{p}(t)$ —plus terms of order $(r/c)^2$ or higher [e.g., $\frac{1}{2}(r/c)^2 \ddot{\mathbf{p}}$] which will be very small for r small enough that \mathbf{p} does not alter markedly in the time r/c .

So Eq. (21.23) gives fields very much like the instantaneous theory—much closer than the instantaneous theory with a delay; the first-order effects of the delay are taken out by the second term. The static formulas are very accurate, much more accurate than you might think. Of course, the compensation only works for points close in. For points far out the correction becomes very bad, because the time delays produce a very large effect, and we get the important $1/r$ term of the radiation.

We still have the problem of computing the electric field and demonstrating that it is the same as Eq. (21.1'). For large distances we can see that the answer is going to come out all right. We know that far from the sources, where we have a propagating wave, \mathbf{E} is perpendicular to \mathbf{B} (and also to \mathbf{r}), as in Fig. 21-4, and that $c\mathbf{B} = \mathbf{E}$. So \mathbf{E} is proportional to the acceleration $\ddot{\mathbf{p}}$, as expected from Eq. (21.1').

To get the electric field completely for all distances, we need to solve for the electrostatic potential. When we computed the current integral for \mathbf{A} to get Eq. (21.18), we made an approximation by disregarding the slight variation of r in the delay terms. This will not work for the electrostatic potential, because we would then get $1/r$ times the integral of the charge density, which is a constant. This approximation is too rough. We need to go to one higher order. Instead of getting involved in that higher-order computation directly, we can do something else—we can determine the scalar potential from Eq. (21.6), using the vector potential we have already found. The divergence of \mathbf{A} , in our case, is just $\partial A_z / \partial z$ —since A_x and A_y are identically zero. Differentiating in the same way that we did above to find \mathbf{B} ,

$$\begin{aligned} \nabla \cdot \mathbf{A} &= \frac{1}{4\pi\epsilon_0 c^2} \left[\dot{\mathbf{p}}(t - r/c) \frac{\partial}{\partial z} \left(\frac{1}{r} \right) + \frac{1}{r} \frac{\partial}{\partial z} \dot{\mathbf{p}}(t - r/c) \right] \\ &= \frac{1}{4\pi\epsilon_0 c^2} \left[-\frac{z\dot{\mathbf{p}}(t - r/c)}{r^3} - \frac{z\ddot{\mathbf{p}}(t - r/c)}{cr^2} \right]. \end{aligned}$$

Or, in vector notation,

$$\nabla \cdot \mathbf{A} = -\frac{1}{4\pi\epsilon_0 c^2} \frac{[\dot{\mathbf{p}} + (r/c)\ddot{\mathbf{p}}]_{t-r/c} \cdot \mathbf{r}}{r^3}.$$

Using Eq. (21.6), we have an equation for ϕ :

$$\frac{\partial \phi}{\partial t} = \frac{1}{4\pi\epsilon_0} \frac{[\dot{\mathbf{p}} + (r/c)\ddot{\mathbf{p}}]_{t-r/c} \cdot \mathbf{r}}{r^3}.$$

Integrating with respect to t just removes one dot from each of the \mathbf{p} 's, so

$$\phi(r, t) = \frac{1}{4\pi\epsilon_0} \frac{[\mathbf{p} + (r/c)\dot{\mathbf{p}}]_{t-r/c} \cdot \mathbf{r}}{r^3} \quad (21.25)$$

(The constant of integration would correspond to some superposed static field which could, of course, exist. For the oscillating dipole we have taken, there is no static field.)

We are now able to find the electric field E from

$$\mathbf{E} = -\nabla\phi - \frac{\partial\mathbf{A}}{\partial t}.$$

Since the steps are tedious but straightforward [providing you remember that $\mathbf{p}(t - r/c)$ and its time derivatives depend on x , y , and z through the retardation r/c], we will just give the result:

$$\mathbf{E}(r, t) = \frac{-1}{4\pi\epsilon_0 r^3} \left[-\mathbf{p}^* - 3 \frac{(\mathbf{p}^* \cdot \mathbf{r})\mathbf{r}}{r^2} + \frac{1}{c^2} \{ \ddot{\mathbf{p}}(t - r/c) \times \mathbf{r} \} \times \mathbf{r} \right] \quad (21.26)$$

with

$$\mathbf{p}^* = \mathbf{p}(t - r/c) + \frac{r}{c} \dot{\mathbf{p}}(t - r/c). \quad (21.27)$$

Although it looks rather complicated, the result is easily interpreted. The vector \mathbf{p}^* is the dipole moment retarded and then “corrected” for the retardation, so the two terms with \mathbf{p}^* give just the static dipole field when r is small. [See Chapter 6, Eq. (6.14).] When r is large, the term in $\ddot{\mathbf{p}}$ dominates, and the electric field is proportional to the acceleration of the charges, at right angles to \mathbf{r} , and, in fact, directed along the projection of $\ddot{\mathbf{p}}$ in a plane perpendicular to \mathbf{r} .

This result agrees with what we would have gotten using Eq. (21.1). Of course, Eq. (21.1) is more general; it works with any motion, while Eq. (21.26) is valid only for small motions for which we can take the retardation r/c as constant over the source. At any rate, we have now provided the underpinnings for our entire previous discussion of light (excepting some matters discussed in Chapter 36 of Vol. I), for it all hinged on the last term of Eq. (21.26). We will discuss next how the fields can be obtained for more rapidly moving charges (leading to the relativistic effects of Chapter 36 of Vol. I).

21-5 The potentials of a moving charge; the general solution of Liénard and Wiechert

In the last section we made a simplification in calculating our integral for A by considering only low velocities. But in doing so we missed an important point and also one where it is easy to go wrong. We will therefore take up now a calculation of the potentials for a point charge moving in any way whatever—even with a relativistic velocity. Once we have this result, we will have the complete electro-magnetism of electric charges. Even Eq. (21.1) can then be derived by taking derivatives. The story will be complete. So bear with us.

Let's try to calculate the scalar potential $\phi(1)$ at the point (x_1, y_1, z_1) produced by a *point* charge, such as an electron, moving in any manner whatsoever. By a “point” charge we mean a very small ball of charge, shrunk down as small as you like, with a charge density $\rho(x, y, z)$. We can find ϕ from Eq. (21.15):

$$\phi(1, t) = \frac{1}{4\pi\epsilon_0} \int \frac{\rho(2, t - r_{12}/c)}{r_{12}} dV_2. \quad (21.28)$$

The answer would seem to be—and almost everyone would, at first, think—that the integral of ρ over such a “point” charge is just the total charge q , so that

$$\phi(1, t) = \frac{1}{4\pi\epsilon_0} \frac{q}{r'_{12}} \quad (\text{wrong}).$$

By r'_{12} we mean the radius vector from the charge at point (2) to point (1) at the retarded time $(t - r_{12}/c)$. It is wrong.

The correct answer is

$$\phi(1, t) = \frac{1}{4\pi\epsilon_0} \frac{q}{r'_{12}} \cdot \frac{1}{1 - v_{r'}/c}, \quad (21.29)$$

where $v_{r'}$ is the component of the velocity of the charge parallel to r'_{12} —namely, toward point (1). We will now show you why. To make the argument easier to

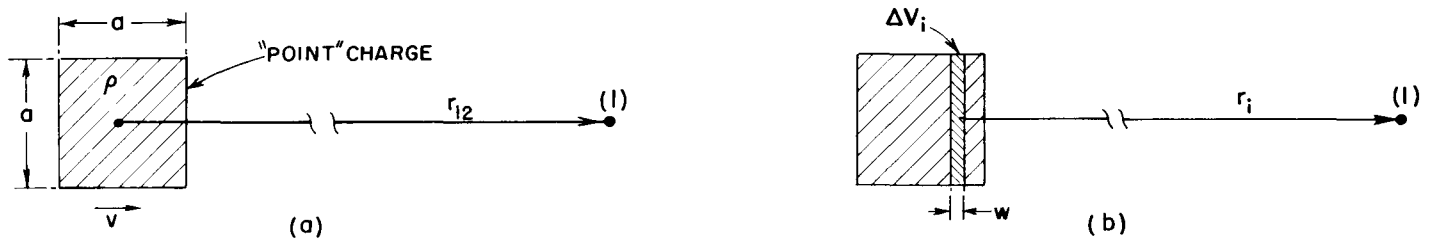


Fig. 21-5. (a) A "point" charge—considered as a small cubical distribution of charge—moving with the speed v toward point (1) (b) The volume element ΔV_i used for calculating the potentials.

follow, we will make the calculation first for a "point" charge which is in the form of a little cube of charge moving toward the point (1) with the speed v , as shown in Fig. 21-5(a). Let the length of a side of the cube be a , which we take to be much, much less than r_{12} , the distance from the center of the charge to the point (1).

Now to evaluate the integral of Eq. (21.28), we will return to basic principles; we will write it as the sum

$$\sum_i \frac{\rho_i \Delta V_i}{r_i}, \quad (21.30)$$

where r_i is the distance from point (1) to the i th volume element ΔV_i , and ρ_i is the charge density at ΔV_i at the time $t_i = t - r_i/c$. Since $r_i \gg a$, always, it will be convenient to take our ΔV_i in the form of thin, rectangular slices perpendicular to r_{12} , as shown in Fig. 21-5(b).

Suppose we start by taking the volume elements ΔV_i with some thickness w much less than a . The individual elements will appear as shown in Fig. 21-6(a), where we have put in more than enough to cover the charge. But we have *not* shown the charge, and for a good reason. Where should we draw it? For each volume element ΔV_i , we are to take ρ at the time $t_i = (t - r_i/c)$, but since the charge is *moving*, it is in a *different place for each volume element* ΔV_i !

Let's say that we begin with the volume element labeled "1" in Fig. 21-6(a), chosen so that at the time $t_1 = (t - r_1/c)$ the "back" edge of the charge occupies ΔV_1 , as shown in Fig. 21-6(b). Then when we evaluate $\rho_2 \Delta V_2$, we must use the position of the charge at the slightly *later* time $t_2 = (t - r_2/c)$, when the charge will be in the position shown in Fig. 21-6(c). And so on, for $\Delta V_3, \Delta V_4$, etc. Now we can evaluate the sum.

Since the thickness of each ΔV_i is w , its volume is wa^2 . Then each volume element that overlaps the charge distribution contains the amount of charge $wa^2\rho$, where ρ is the density of charge within the cube—which we take to be uniform. When the distance from the charge to point (1) is large, we will make a negligible error by setting all the r_i 's in the denominators equal to some average value, say the retarded position r' of the center of the charge. Then the sum (21.30) is

$$\sum_{i=1}^N \frac{\rho wa^2}{r'},$$

where ΔV_N is the last ΔV_i that overlaps the charge distributions, as shown in Fig. 21-6(e). The sum is, clearly,

$$N \frac{\rho wa^2}{r'} = \frac{\rho a^3}{r'} \left(\frac{Nw}{a} \right).$$

Now ρa^3 is just the total charge q and Nw is the length b shown in part (e) of the figure. So we have

$$\phi = \frac{q}{4\pi\epsilon_0 r'} \left(\frac{b}{a} \right). \quad (21.31)$$

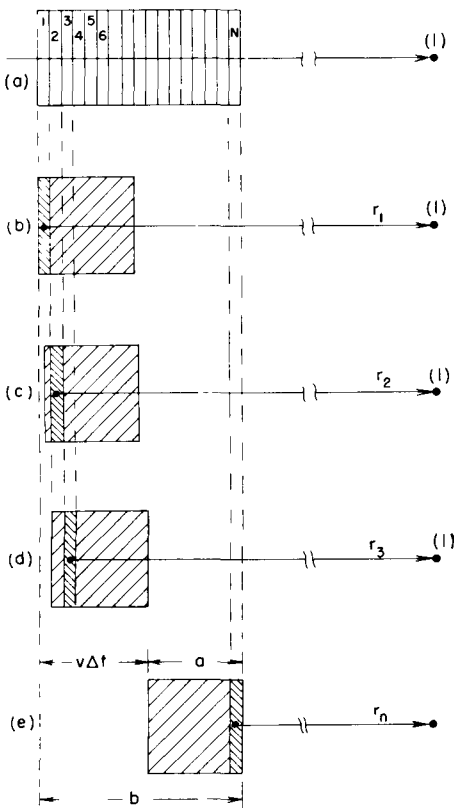


Fig. 21-6. Integrating $\rho(t - r'/c) dV$ for a moving charge.

What is b ? It is the length of the cube of charge *increased* by the distance moved by the charge between $t_1 = (t - r_1/c)$ and $t_N = (t - r_N/c)$ —which is the distance the charge moves in the time

$$\Delta t = t_N - t_1 = (r_1 - r_N)/c = b/c.$$

Since the speed of the charge is v , the distance moved is $v \Delta t = vb/c$. But the length b is this distance added to a :

$$b = a + \frac{v}{c} b.$$

Solving for b , we get

$$b = \frac{a}{1 - (v/c)}.$$

Of course by v we mean the velocity at the retarded time $t' = (t - r'/c)$, which we can indicate by writing $[1 - v/c]_{\text{ret}}$, and Eq. (21.31) for the potential becomes

$$\phi(1, t) = \frac{q}{4\pi\epsilon_0 r'} \frac{1}{[1 - (v/c)]_{\text{ret}}}.$$

This result agrees with our assertion, Eq. (21.29). There is a correction term which comes about because the charge is moving as our integral “sweeps over the charge.” When the charge is moving toward the point (1), its contribution to the integral is increased by the ratio b/a . Therefore the correct integral is q/r' multiplied by b/a , which is $1/[1 - v/c]_{\text{ret}}$.

If the velocity of the charge is not directed toward the observation point (1), you can see that what matters is the *component* of its velocity toward point (1). Calling this velocity component v_r , the correction factor is $1/[1 - v_r/c]_{\text{ret}}$. Also, the analysis we have made goes exactly the same way for a charge distribution of *any* shape—it doesn't have to be a cube. Finally, since the “size” of the charge a doesn't enter into the final result, the same result holds when we let the charge shrink to any size—even to a point. The general result is that the scalar potential for a point charge moving with any velocity is

$$\phi(t) = \frac{q}{4\pi\epsilon_0 r' [1 - (v_r/c)]_{\text{ret}}}. \quad (21.32)$$

This equation is often written in the equivalent form

$$\phi(1, t) = \frac{q}{4\pi\epsilon_0 [r - (\mathbf{v} \cdot \mathbf{r}/c)]_{\text{ret}}}, \quad (21.33)$$

where \mathbf{r} is the vector from the charge to the point (1), where ϕ is being evaluated, and all the quantities in the bracket are to have their values at the retarded time $t' = t - r'/c$.

The same thing happens when we compute A for a point charge, from Eq. (21.16). The current density is $\rho\mathbf{v}$ and the integral over ρ is the same as we found for ϕ . The vector potential is

$$A(1, t) = \frac{qv}{4\pi\epsilon_0 c^2 [r - (\mathbf{v} \cdot \mathbf{r}/c)]_{\text{ret}}}. \quad (21.34)$$

The potentials for a point charge were first deduced in this form by Liénard and Wiechert and are called the *Liénard-Wiechert potentials*.

To close the ring back to Eq. (21.1) it is only necessary to compute \mathbf{E} and \mathbf{B} from these potentials (using $\mathbf{B} = \nabla \times \mathbf{A}$ and $\mathbf{E} = -\nabla\phi - \partial\mathbf{A}/\partial t$). It is now only arithmetic. The arithmetic, however, is fairly involved, so we will not write out the details. Perhaps you will take our word for it that Eq. (21.1) is equivalent to the Liénard-Wiechert potentials we have derived.*

* If you have a lot of paper and time you can try to work it through yourself. We would, then, make two suggestions. First, don't forget that the derivatives of r' are complicated, since it is a function of t' . Second, don't try to *derive* (21.1), but carry out all of the derivatives in it, and then compare what you get with the \mathbf{E} obtained from the potentials (21.33) and (21.34).

21-6 The potentials for a charge moving with constant velocity; the Lorentz formula

We want next to use the Liénard-Wiechert potentials for a special case—to find the fields of a charge moving with uniform velocity in a straight line. We will do it again later, using the principle of relativity. We already know what the potentials are when we are standing in the rest frame of a charge. When the charge is moving, we can figure everything out by a relativistic transformation from one system to the other. But relativity had its origin in the theory of electricity and magnetism. The formulas of the Lorentz transformation (Chapter 15, Vol. I) were discoveries made by Lorentz when he was studying the equations of electricity and magnetism. So that you can appreciate where things have come from, we would like to show that the Maxwell equations do lead to the Lorentz transformation. We begin by calculating the potentials of a charge moving with uniform velocity, directly from the electrodynamics of Maxwell's equations. We have shown that Maxwell's equations lead to the potentials for a moving charge that we got in the last section. So when we use these potentials, we are using Maxwell's theory.

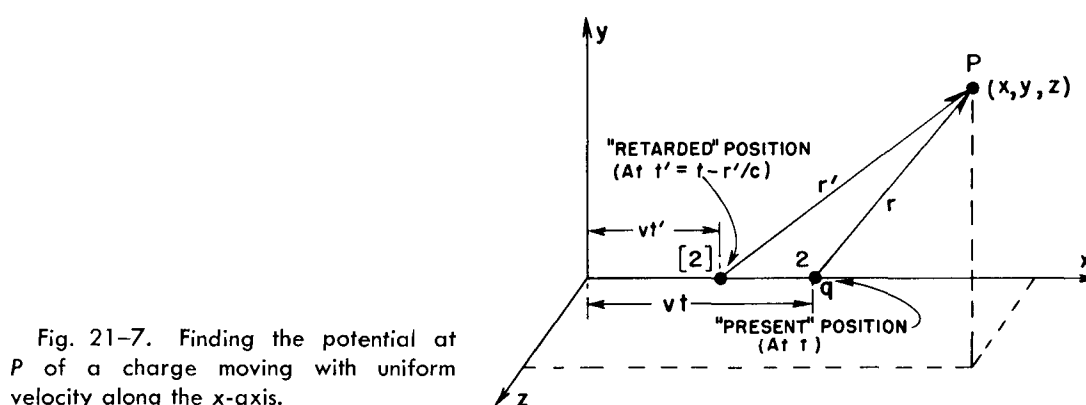


Fig. 21-7. Finding the potential at P of a charge moving with uniform velocity along the x -axis.

Suppose we have a charge moving along the x -axis with the speed v . We want the potentials at the point $P(x, y, z)$, as shown in Fig. 21-7. If $t = 0$ is the moment when the charge is at the origin, at the time t the charge is at $x = vt, y = z = 0$. What we need to know, however, is its position at the retarded time

$$t' = t - \frac{r'}{c}, \quad (21.35)$$

where r' is the distance to the point P from the charge *at the retarded time*. At the earlier time t' , the charge was at $x = vt'$, so

$$r' = \sqrt{(x - vt')^2 + y^2 + z^2}. \quad (21.36)$$

To find r' or t' we have to combine this equation with Eq. (21.35). First, we eliminate r' by solving Eq. (21.35) for r' and substituting in Eq. (21.36). Then, squaring both sides, we get

$$c^2(t - t')^2 = (x - vt')^2 + y^2 + z^2,$$

which is a quadratic equation in t' . Expanding the squared binomials and collecting like terms in t' , we get

$$(c^2 - v^2)t'^2 - 2(xv - c^2t)t' + x^2 + y^2 + z^2 - (ct)^2 = 0.$$

Solving for t' ,

$$\left(1 - \frac{v^2}{c^2}\right)t' = t - \frac{vx}{c^2} - \frac{1}{c} \sqrt{(x - vt)^2 + \left(1 - \frac{v^2}{c^2}\right)(y^2 + z^2)} \quad (21.37)$$

To get r' we have to substitute this expression for t' into

$$r' = c(t - t').$$

Now we are ready to find ϕ from Eq. (21.33), which, since v is constant, becomes

$$\phi(x, y, z, t) = \frac{q}{4\pi\epsilon_0} \frac{1}{r' - (\mathbf{v} \cdot \mathbf{r}'/c)}. \quad (21.38)$$

The component of \mathbf{v} in the direction of \mathbf{r}' is $v \times (x - vt)/r'$, so $\mathbf{v} \cdot \mathbf{r}'$ is just $v \times (x - vt)$, and the whole denominator is

$$c(t - t') - \frac{v}{c}(x - vt) = c \left[t - \frac{vx}{c^2} - \left(1 - \frac{v^2}{c^2}\right) t' \right].$$

Substituting for $(1 - v^2/c^2)t'$ from Eq. (21.37), we get for ϕ

$$\phi(x, y, z, t) = \frac{q}{4\pi\epsilon_0} \frac{1}{\sqrt{(x - vt)^2 + \left(1 - \frac{v^2}{c^2}\right)(y^2 + z^2)}}.$$

This equation is more understandable if we rewrite it as

$$\phi(x, y, z, t) = \frac{q}{4\pi\epsilon_0} \frac{1}{\sqrt{1 - \frac{v^2}{c^2}}} \frac{1}{\left[\left(\frac{x - vt}{\sqrt{1 - v^2/c^2}} \right)^2 + y^2 + z^2 \right]^{1/2}}. \quad (21.39)$$

The vector potential \mathbf{A} is the same expression with an additional factor of v/c^2 .

$$\mathbf{A} = \frac{v}{c^2} \phi.$$

In Eq. (21.39) you can clearly see the beginning of the Lorentz transformation. If the charge were at the origin in its own rest frame, its potential would be

$$\phi(x, y, z) = \frac{q}{4\pi\epsilon_0} \frac{1}{[x^2 + y^2 + z^2]^{1/2}}.$$

We are seeing it in a moving coordinate system, and it appears that the coordinates should be transformed by

$$\begin{aligned} x &\rightarrow \frac{x - vt}{\sqrt{1 - v^2/c^2}}, \\ y &\rightarrow y, \\ z &\rightarrow z. \end{aligned}$$

That is just the Lorentz transformation, and what we have done is essentially the way Lorentz discovered it.

But what about that extra factor $1/\sqrt{1 - v^2/c^2}$ that appears at the front of Eq. (21.39)? Also, how does the vector potential \mathbf{A} appear, when it is everywhere zero in the rest frame of the particle? We will soon show that \mathbf{A} and ϕ *together* constitute a four-vector, like the momentum \mathbf{p} and the total energy U of a particle. The extra $1/\sqrt{1 - v^2/c^2}$ in Eq. (21.39) is the same factor that always comes in when one transforms the components of a four-vector—just as the charge density ρ transforms to $\rho/\sqrt{1 - v^2/c^2}$. In fact, it is almost apparent from Eqs. (21.4) and (21.5) that \mathbf{A} and ϕ are components of a four-vector, because we have already shown in Chapter 13 that \mathbf{j} and ρ are the components of a four-vector.

Later we will take up in more detail the relativity of electrodynamics; here we only wished to show how naturally the Maxwell equations lead to the Lorentz transformation. You will not, then, be surprised to find that the laws of electricity and magnetism are already correct for Einstein's relativity. We will not have to "fix them up," as we had to do for Newton's laws of mechanics.

AC Circuits

22-1 Impedances

Most of our work in this course has been aimed at reaching the complete equations of Maxwell. In the last two chapters we have been discussing the consequences of these equations. We have found that the equations contain all the static phenomena we had worked out earlier, as well as the phenomena of electromagnetic waves and light that we had gone over in some detail in Volume I. The Maxwell equations give both phenomena, depending upon whether one computes the fields close to the currents and charges, or very far from them. There is not much interesting to say about the intermediate region; no special phenomena appear there.

There still remain, however, several subjects in electromagnetism that we want to take up. We want to discuss the question of relativity and the Maxwell equations—what happens when one looks at the Maxwell equations with respect to moving coordinate systems. There is also the question of the conservation of energy in electromagnetic systems. Then there is the broad subject of the electromagnetic properties of materials; so far, except for the study of the properties of dielectrics, we have considered only the electromagnetic fields in free space. And although we covered the subject of light in some detail in Volume I, there are still a few things we would like to do again from the point of view of the field equations.

In particular, we want to take up again the subject of the index of refraction, particularly for dense materials. Finally, there are the phenomena associated with waves confined in a limited region of space. We touched on this kind of problem briefly when we were studying sound waves. Maxwell's equations lead also to solutions which represent confined waves of the electric and magnetic fields. We will take up this subject, which has important technical applications, in some of the following chapters. In order to lead up to that subject, we will begin by considering the properties of electrical circuits at low frequencies. We will then be able to make a comparison between those situations in which the almost static approximations of Maxwell's equations are applicable and those situations in which high-frequency effects are dominant.

So we descend from the great and esoteric heights of the last few chapters and turn to the relatively low-level subject of electrical circuits. We will see, however, that even such a mundane subject, when looked at in sufficient detail, can contain great complications.

We have already discussed some of the properties of electrical circuits in Chapters 23 and 25 of Vol. I. Now we will cover some of the same material again, but in greater detail. Again we are going to deal only with linear systems and with voltages and currents which all vary sinusoidally; we can then represent all voltages and currents by complex numbers, using the exponential notation described in Chapter 22 of Vol. I. Thus a time-varying voltage $V(t)$ will be written

$$V(t) = \hat{V}e^{i\omega t}, \quad (22.1)$$

where \hat{V} represents a complex number that is independent of t . It is, of course, understood that the actual time-varying voltage $V(t)$ is given by the *real part* of the complex function on the right-hand side of the equation.

22-1 Impedances

22-2 Generators

22-3 Networks of ideal elements; Kirchhoff's rules

22-4 Equivalent circuits

22-5 Energy

22-6 A ladder network

22-7 Filters

22-8 Other circuit elements

Review: Chapter 22, Vol. I, *Algebra*
Chapter 23, Vol. I, *Resonance*
Chapter 25, Vol. I, *Linear*
Systems and Review

Similarly, all of our other time-varying quantities will be taken to vary sinusoidally at the same frequency ω . So we write

$$\begin{aligned} I &= \hat{I} e^{i\omega t} \quad (\text{current}), \\ \mathcal{E} &= \hat{\mathcal{E}} e^{i\omega t} \quad (\text{emf}), \\ \mathbf{E} &= \hat{\mathbf{E}} e^{i\omega t} \quad (\text{electric field}), \end{aligned} \quad (22.2)$$

and so on.

Most of the time we will write our equations in terms of V, I, \mathcal{E}, \dots (instead of in terms of $\hat{V}, \hat{I}, \hat{\mathcal{E}}, \dots$), remembering, though, that the time variations are as given in (22.2).

In our earlier discussion of circuits we assumed that such things as inductances, capacitances, and resistances were familiar to you. We want now to look in a little more detail at what is meant by these idealized circuit elements. We begin with the inductance.

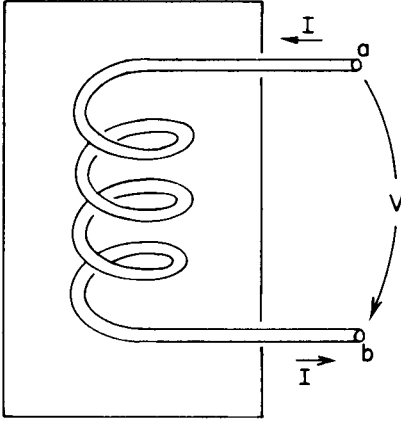


Fig. 22-1. An inductance.

An inductance is made by winding many turns of wire in the form of a coil and bringing the two ends out to terminals at some distance from the coil, as shown in Fig. 22-1. We want to assume that the magnetic field produced by currents in the coil does not spread out strongly all over space and interact with other parts of the circuit. This is usually arranged by winding the coil in a doughnut-shaped form, or by confining the magnetic field by winding the coil on a suitable iron core, or by placing the coil in some suitable metal box, as indicated schematically in Fig. 22-1. In any case, we assume that there is a negligible magnetic field in the external region near the terminals a and b . We are also going to assume that we can neglect any electrical resistance in the wire of the coil. Finally, we will assume that we can neglect the amount of electrical charge that appears on the surface of a wire in building up the electric fields.

With all these approximations we have what we call an “ideal” inductance. (We will come back later and discuss what happens in a real inductance.) For an ideal inductance we say that the voltage across the terminals is equal to $L(dI/dt)$. Let’s see why that is so. When there is a current through the inductance, a magnetic field proportional to the current is built up inside the coil. If the current changes with time, the magnetic field also changes. In general, the curl of \mathbf{E} is equal to $-d\mathbf{B}/dt$; or, put differently, the line integral of \mathbf{E} all the way around any closed path is equal to the negative of the rate of change of the flux of \mathbf{B} through the loop. Now suppose we consider the following path: Begin at terminal a and go along the coil (staying always inside the wire) to terminal b ; then return from terminal b to terminal a through the air in the space outside the inductance. The line integral of \mathbf{E} around this closed path can be written as the sum of two parts:

$$\int \mathbf{E} \cdot d\mathbf{s} = \int_a^b \underset{\text{via coil}}{\mathbf{E} \cdot d\mathbf{s}} + \int_b^a \underset{\text{outside}}{\mathbf{E} \cdot d\mathbf{s}}. \quad (22.3)$$

As we have seen before, there can be no electric fields inside a perfect conductor. (The smallest fields would produce infinite currents.) Therefore the integral from a to b via the coil is zero. The whole contribution to the line integral of \mathbf{E} comes from the path outside the inductance from terminal b to terminal a . Since we have assumed that there are no magnetic fields in the space outside of the “box,” this part of the integral is independent of the path chosen and we can define the potentials of the two terminals. The difference of these two potentials is what we call the voltage difference, or simply the voltage V , so we have

$$V = - \int_b^a \mathbf{E} \cdot d\mathbf{s} = - \oint \mathbf{E} \cdot d\mathbf{s}.$$

The complete line integral is what we have before called the electromotive force \mathcal{E} and is, of course, equal to the rate of change of the magnetic flux in the coil. We have seen earlier that this emf is equal to the negative rate of change of

the current, so we have

$$V = -\varepsilon = L \frac{dI}{dt},$$

where L is the inductance of the coil. Since $dI/dt = i\omega I$, we have

$$V = i\omega LI. \quad (22.4)$$

The way we have described the ideal inductance illustrates the general approach to other ideal circuit elements—usually called “lumped” elements. The properties of the element are described completely in terms of currents and voltages that appear at the terminals. By making suitable approximations, it is possible to ignore the great complexities of the fields that appear inside the object. A separation is made between what happens inside and what happens outside.

For all the circuit elements we will find a relation like the one in Eq. (22.4), in which the voltage is proportional to the current with a proportionality constant that is, in general, a complex number. This complex coefficient of proportionality is called the *impedance* and is usually written as z (not to be confused with the z -coordinate). It is, in general, a function of the frequency ω . So for any lumped element we write

$$\frac{V}{I} = \frac{\hat{V}}{\hat{I}} = z. \quad (22.5)$$

For an inductance, we have

$$z(\text{inductance}) = z_L = i\omega L. \quad (22.6)$$

Now let's look at a capacitor from the same point of view.* A capacitor consists of a pair of conducting plates from which two wires are brought out to suitable terminals. The plates may be of any shape whatsoever, and are often separated by some dielectric material. We illustrate such a situation schematically in Fig. 22-2. Again we make several simplifying assumptions. We assume that the plates and the wires are perfect conductors. We also assume that the insulation between the plates is perfect, so that no charges can flow across the insulation from one plate to the other. Next, we assume that the two conductors are close to each other but far from all others, so that all field lines which leave one plate end up on the other. Then there are always equal and opposite charges on the two plates and the charges on the plates are much larger than the charges on the surfaces of the lead-in wires. Finally, we assume that there are no magnetic fields close to the capacitor.

Suppose now we consider the line integral of \mathbf{E} around a closed loop which starts at terminal a , goes along inside the wire to the top plate of the capacitor, jumps across the space between the plates, passes from the lower plate to terminal b through the wire, and returns to terminal a in the space outside the capacitor. Since there is no magnetic field, the line integral of \mathbf{E} around this closed path is zero. The integral can be broken down into three parts:

$$\oint \mathbf{E} \cdot d\mathbf{s} = \int_{\text{along wires}} \mathbf{E} \cdot d\mathbf{s} + \int_{\text{between plates}} \mathbf{E} \cdot d\mathbf{s} + \int_b^a \mathbf{E} \cdot d\mathbf{s}. \quad (22.7)$$

The integral along the wires is zero, because there are no electric fields inside perfect conductors. The integral from b to a outside the capacitor is equal to the negative of the potential difference between the terminals. Since we imagined that the two plates are in some way isolated from the rest of the world, the total charge on

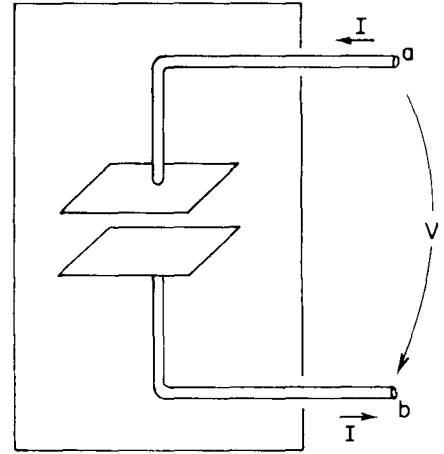


Fig. 22-2. A capacitor (or condenser).

* There are people who say we should call the *objects* by the names “inductor” and “capacitor” and call their *properties* “inductance” and “capacitance” (by analogy with “resistor” and “resistance”). We would rather use the words you will hear in the laboratory. Most people still say “inductance” for both the physical coil and its inductance L . The word “capacitor” seems to have caught on—although you will still hear “condenser” fairly often—and most people still prefer the sound of “capacity” to “capacitance.”

the two plates must be zero; if there is a charge Q on the upper plate, there is an equal, opposite charge $-Q$ on the lower plate. We have seen earlier that if two conductors have equal and opposite charges, plus and minus Q , the potential difference between the plates is equal to Q/C , where C is called the capacity of the two conductors. From Eq. (22.7) the potential difference between the terminals a and b is equal to the potential difference between the plates. We have, therefore, that

$$V = \frac{Q}{C}.$$

The electric current I entering the capacitor through terminal a (and leaving through terminal b) is equal to dQ/dt , the rate of change of the electric charge on the plates. Writing dV/dt as $i\omega V$, we can put the voltage current relationship for a capacitor in the following way:

$$i\omega V = \frac{I}{C},$$

or

$$V = \frac{I}{i\omega C}. \quad (22.8)$$

The impedance z of a capacitor, is then

$$z \text{ (capacitor)} = z_C = \frac{1}{i\omega C}. \quad (22.9)$$

The third element we want to consider is a resistor. However, since we have not yet discussed the electrical properties of real materials, we are not yet ready to talk about what happens inside a real conductor. We will just have to accept as fact that electric fields can exist inside real materials, that these electric fields give rise to a flow of electric charge—that is, to a current—and that this current is proportional to the integral of the electric field from one end of the conductor to the other. We then imagine an ideal resistor constructed as in the diagram of Fig. 22-3. Two wires which we take to be perfect conductors go from the terminals a and b to the two ends of a bar of resistive material. Following our usual line of argument, the potential difference between the terminals a and b is equal to the line integral of the external electric field, which is also equal to the line integral of the electric field through the bar of resistive material. It then follows that the current I through the resistor is proportional to the terminal voltage V :

$$I = \frac{V}{R},$$

where R is called the resistance. We will see later that the relation between the current and the voltage for real conducting materials is only approximately linear. We will also see that this approximate proportionality is expected to be independent of the frequency of variation of the current and voltage only if the frequency is not too high. For alternating currents then, the voltage across a resistor is in phase with the current, which means that the impedance is a real number.

$$z \text{ (resistance)} = z_R = R. \quad (22.10)$$

Our results for the three lumped circuit elements—the inductor, the capacitor, and the resistor—are summarized in Fig. 22-4. In this figure, as well as in the preceding ones, we have indicated the voltage by an arrow that is directed from one terminal to another. If the voltage is “positive”—that is, if the terminal a is at a *higher* potential than the terminal b —the arrow indicates the direction of a positive “voltage drop.”

Although we are talking about alternating currents, we can of course include the special case of circuits with steady currents by taking the limit as the frequency ω goes to zero. For zero frequency—that is, for DC—the impedance of an inductance goes to zero; it becomes a short circuit. For DC, the impedance of a condenser

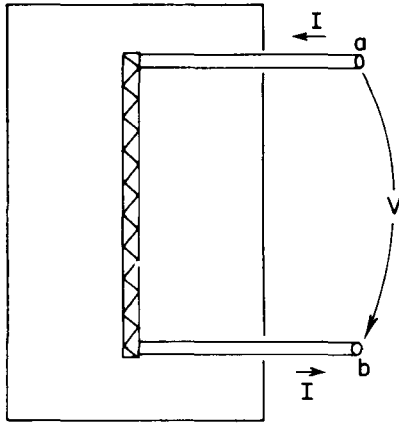


Fig. 22-3. A resistor.

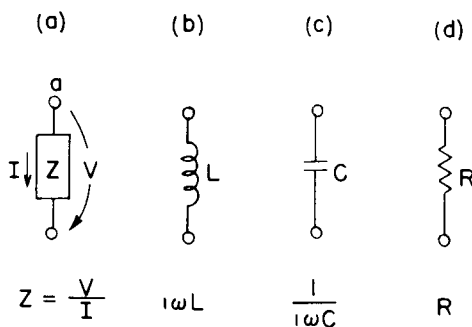


Fig. 22-4. The ideal lumped circuit elements (passive).

goes to infinity; it becomes an open circuit. Since the impedance of a resistor is independent of frequency, it is the only element left when we analyze a circuit for DC.

In the circuit elements we have described so far, the current and voltage are proportional to each other. If one is zero, so also is the other. We usually think in terms like these: An applied voltage is “responsible” for the current, or a current “gives rise to” a voltage across the terminals; so in a sense the elements “respond” to the “applied” external conditions. For this reason these elements are called *passive elements*. They can thus be contrasted with the active elements, such as the generators we will consider in the next section, which are the *sources* of the oscillating currents or voltages in a circuit.

22-2 Generators

Now we want to talk about an *active* circuit element—one that is a source of the currents and voltages in a circuit—namely, a *generator*.

Suppose that we have a coil like an inductance except that it has very few turns, so that we may neglect the magnetic field of its own current. This coil, however, sits in a changing magnetic field such as might be produced by a rotating magnet, as sketched in Fig. 22-5. (We have seen earlier that such a rotating magnetic field can also be produced by a suitable set of coils with alternating currents.) Again we must make several simplifying assumptions. The assumptions we will make are all the ones that we described for the case of the inductance. In particular, we assume that the varying magnetic field is restricted to a definite region in the vicinity of the coil and does not appear outside the generator in the space between the terminals.

Following closely the analysis we made for the inductance, we consider the line integral of \mathbf{E} around a complete loop that starts at terminal a , goes through the coil to terminal b and returns to its starting point in the space between the two terminals. Again we conclude that the potential difference between the terminals is equal to the total line integral of \mathbf{E} around the loop:

$$V = - \oint \mathbf{E} \cdot d\mathbf{s}.$$

This line integral is equal to the emf in the circuit, so the potential difference V across the terminals of the generator is also equal to the rate of change of the magnetic flux linking the coil:

$$V = -\varepsilon = \frac{d}{dt} (\text{flux}). \quad (22.11)$$

For an ideal generator we assume that the magnetic flux linking the coil is determined by external conditions—such as the angular velocity of a rotating magnetic field—and is not influenced in any way by the currents through the generator. Thus a generator—at least the *ideal* generator we are considering—is not an impedance. The potential difference across its terminals is determined by the arbitrarily assigned electromotive force $\varepsilon(t)$. Such an ideal generator is represented by the symbol shown in Fig. 22-6. The little arrow represents the direction of the emf when it is positive. A positive emf in the generator of Fig. 22-6 will produce a voltage $V = \varepsilon$, with the terminal a at a higher potential than the terminal b .

There is another way to make a generator which is quite different on the inside but which is indistinguishable from the one we have just described insofar as what happens beyond its terminals. Suppose we have a coil of wire which is rotated in a *fixed* magnetic field, as indicated in Fig. 22-7. We show a bar magnet to indicate the presence of a magnetic field; it could, of course, be replaced by any other source of a steady magnetic field, such as an additional coil carrying a steady current. As shown in the figure, connections from the rotating coil are made to the outside world by means of sliding contacts or “slip rings.” Again, we are interested in the potential difference that appears across the two terminals

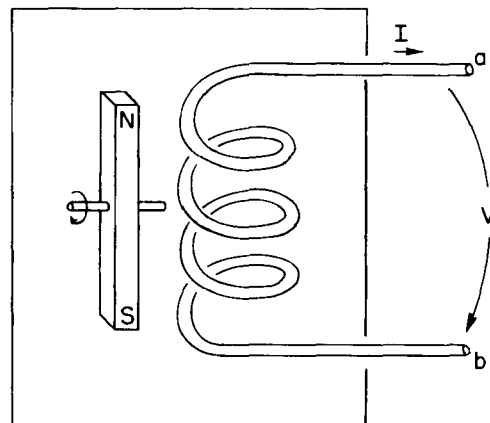


Fig. 22-5. A generator consisting of a fixed coil and a rotating magnetic field.

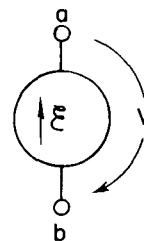


Fig. 22-6. Symbol for an ideal generator.

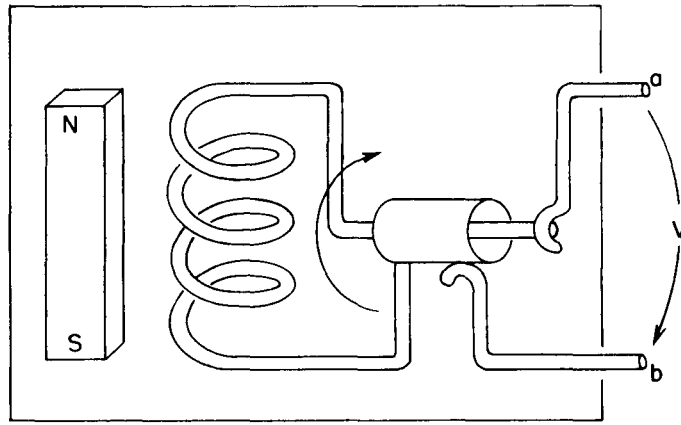


Fig. 22-7. A generator consisting of a coil rotating in a fixed magnetic field.

a and b , which is of course the integral of the electric field from terminal a to terminal b along a path outside the generator.

Now in the system of Fig. 22-7 there are no changing magnetic fields, so we might at first wonder how any voltage could appear at the generator terminals. In fact, there are no electric fields anywhere inside the generator. We are, as usual, assuming for our ideal elements that the wires inside are made of a perfectly conducting material, and as we have said many times, the electric field inside a perfect conductor is equal to zero. But that is not true. It is not true when a conductor is moving in a magnetic field. The true statement is that the total *force* on any charge inside a perfect conductor must be zero. Otherwise there would be an infinite flow of the free charges. So what is always true is that the sum of the electric field E and the cross product of the velocity of the conductor and the magnetic field B —which is the total force on a unit charge—must have the value zero inside the conductor:

$$\mathbf{F} = \mathbf{E} + \mathbf{v} \times \mathbf{B} = 0 \quad (\text{in a perfect conductor}), \quad (22.12)$$

where \mathbf{v} represents the velocity of the conductor. Our earlier statement that there is no electric field inside a perfect conductor is all right if the velocity \mathbf{v} of the conductor is zero; otherwise the correct statement is given by Eq. (22.12).

Returning to our generator of Fig. 22-7, we now see that the line integral of the electric field E from terminal a to terminal b through the conducting path of the generator must be equal to the line integral of $\mathbf{v} \times \mathbf{B}$ on the same path,

$$\int_a^b \underset{\text{inside conductor}}{E \cdot ds} = - \int_a^b \underset{\text{inside conductor}}{(\mathbf{v} \times \mathbf{B}) \cdot ds}. \quad (22.13)$$

It is still true, however, that the line integral of E around a complete loop, including the return from b to a outside the generator, must be zero, because there are no changing magnetic fields. So the first integral in Eq. (22.13) is also equal to V , the voltage between the two terminals. It turns out that the right-hand integral of Eq. (22.13) is just the rate of change of the flux linkage through the coil and is therefore—by the flux rule—equal to the emf in the coil. So we have again that the potential difference across the terminals is equal to the electromotive force in the circuit, in agreement with Eq. (22.11). So whether we have a generator in which a magnetic field changes near a fixed coil, or one in which a coil moves in a fixed magnetic field, the external properties of the generators are the same. There is a voltage difference V across the terminals, which is independent of the current in the circuit but depends only on the arbitrarily assigned conditions inside the generator.

So long as we are trying to understand the operation of generators from the point of view of Maxwell's equations, we might also ask about the ordinary chemical cell, like a flashlight battery. It is also a generator, i.e., a voltage source, although it will of course only appear in DC circuits. The simplest kind of cell to understand is shown in Fig. 22-8. We imagine two metal plates immersed in some

chemical solution. We suppose that the solution contains positive and negative ions. We suppose also that one kind of ion, say the negative, is much heavier than the one of opposite polarity, so that its motion through the solution by the process of diffusion is much slower. We suppose next that by some means or other it is arranged that the concentration of the solution is made to vary from one part of the liquid to the other, so that the number of ions of both polarities near, say, the lower plate is much larger than the concentration of ions near the upper plate. Because of their rapid mobility the positive ions will drift more readily into the region of lower concentration, so that there will be a slight excess of positive charge arriving at the upper plate. The upper plate will become positively charged and the lower plate will have a net negative charge.

As more and more charges diffuse to the upper plate, the potential of this plate will rise until the resulting electric field between the plates produces forces on the ions which just compensate for their excess mobility, so the two plates of the cell quickly reach a potential difference which is characteristic of the internal construction.

Arguing just as we did for the ideal capacitor, we see that the potential difference between the terminals *a* and *b* is just equal to the line integral of the electric field between the two plates when there is no longer any net diffusion of the ions. There is, of course, an essential difference between a capacitor and such a chemical cell. If we short-circuit the terminals of a condenser for a moment, the capacitor is discharged and there is no longer any potential difference across the terminals. In the case of the chemical cell a current can be drawn from the terminals continuously without any change in the emf—until, of course, the chemicals inside the cell have been used up. In a real cell it is found that the potential difference across the terminals decreases as the current drawn from the cell increases. In keeping with the abstractions we have been making, however, we may imagine an ideal cell in which the voltage across the terminals is independent of the current. A real cell can then be looked at as an ideal cell in series with a resistor.

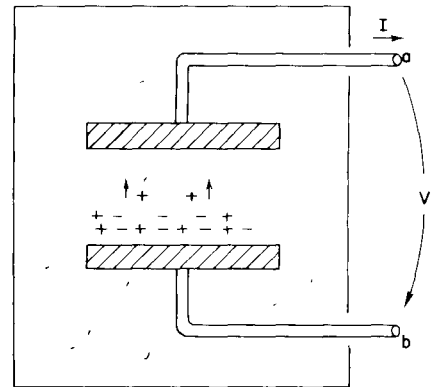


Fig. 22-8. A chemical cell.

22-3 Networks of ideal elements; Kirchhoff's rules

As we have seen in the last section, the description of an ideal circuit element in terms of what happens outside the element is quite simple. The current and the voltage are linearly related. But what is actually happening inside the element is quite complicated, and it is quite difficult to give a precise description in terms of Maxwell's equations. Imagine trying to give a precise description of the electric and magnetic fields of the inside of a radio which contains hundreds of resistors, capacitors, and inductors. It would be an impossible task to analyze such a thing by using Maxwell's equations. But by making the many approximations we have described in Section 22-2 and summarizing the essential features of the real circuit elements in terms of idealizations, it becomes possible to analyze an electrical circuit in a relatively straightforward way. We will now show how that is done.

Suppose we have a circuit consisting of a generator and several impedances connected together, as shown in Fig. 22-9. According to our approximations there is no magnetic field in the region outside the individual circuit elements. Therefore the line integral of \mathbf{E} around any curve which does not pass through any of the elements is zero. Consider then the curve Γ shown by the broken line which goes all the way around the circuit in Fig. 22-9. The line integral of \mathbf{E} around this curve is made up of several pieces. Each piece is the line integral from one terminal of a circuit element to the other. This line integral we have called the voltage drop across the circuit element. The complete line integral is then just the sum of the voltage drops across all of the elements in the circuit:

$$\oint \mathbf{E} \cdot d\mathbf{s} = \sum V_n.$$

Since the line integral is zero, we have that the sum of the potential differences

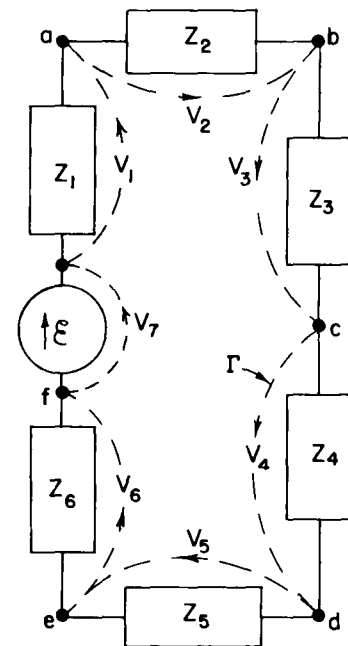


Fig. 22-9. The sum of the voltage drops around any closed path is zero.

around a complete loop of a circuit is equal to zero:

$$\sum_{\text{around any loop}} V_n = 0. \quad (22.14)$$

This result follows from one of Maxwell's equations—that in a region where there are no magnetic fields the line integral of \mathbf{E} around any complete loop is zero.

Suppose we consider now a circuit like that shown in Fig. 22-10. The horizontal line joining the terminals $a, b, c,$ and d is intended to show that these terminals are all connected, or that they are joined by wires of negligible resistance. In any case, the drawing means that terminals $a, b, c,$ and d are all at the same potential and, similarly, that the terminals $e, f, g,$ and h are also at one common potential. Then the voltage drop V across each of the four elements is the same.

Now one of our idealizations has been that negligible electrical charges accumulate on the terminals of the impedances. We now assume further that any electrical charges on the wires joining terminals can also be neglected. Then the conservation of charge requires that any charge which leaves one circuit element immediately enters some other circuit element. Or, what is the same thing, we require that the algebraic sum of the currents which enter any given junction must be zero. By a junction, of course, we mean any set of terminals such as $a, b, c,$ and d which are connected. Such a set of connected terminals is usually called a "node." The conservation of charge then requires that for the circuit of Fig. 22-10,

$$I_1 - I_2 - I_3 - I_4 = 0. \quad (22.15)$$

The sum of the currents entering the node which consists of the four terminals $e, f, g,$ and h must also be zero:

$$-I_1 + I_2 + I_3 + I_4 = 0. \quad (22.16)$$

This is, of course, the same as Eq. (22.15). The two equations are not independent. The general rule is that *the sum of the currents into any node must be zero*.

$$\sum_{\text{into a node}} I_n = 0. \quad (22.17)$$

Our earlier conclusion that the sum of the voltage drops around a closed loop is zero must apply to any loop in a complicated circuit. Also, our result that the sum of the currents into a node is zero must be true for any node. These two equations are known as *Kirchhoff's rules*. With these two rules it is possible to solve for the currents and voltages in any network whatever.

Suppose we consider the more complicated circuit of Fig. 22-11. How shall we find the currents and voltages in this circuit? We can find them in the following straightforward way. We consider separately each of the four subsidiary closed loops which appear in the circuit. (For instance, one loop goes from terminal a to terminal b to terminal e to terminal d and back to terminal a .) For each of the loops we write the equation for the first of Kirchhoff's rules—that the sum of the voltages around each loop is equal to zero. We must remember to count the voltage drop as positive if we are going *in* the direction of the current and negative if we are going across an element in the direction *opposite* to the current; and we must remember that the voltage drop across a generator is the *negative* of the emf in that direction. Thus if we consider the small loop that starts and ends at terminal a we have the equation

$$z_1 I_1 + z_3 I_3 + z_4 I_4 - \mathcal{E}_1 = 0.$$

Applying the same rule to the remaining loops, we would get three more equations of the same kind.

Next, we must write the current equation for each of the nodes in the circuit. For example, summing the currents into the node at terminal b gives the equation

$$I_1 - I_3 - I_2 = 0.$$

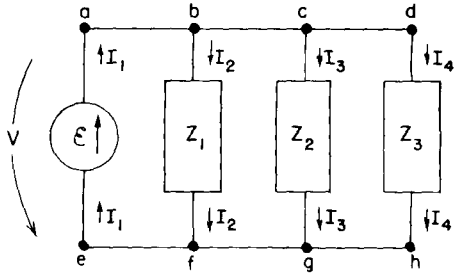


Fig. 22-10. The sum of the currents into any node is zero.

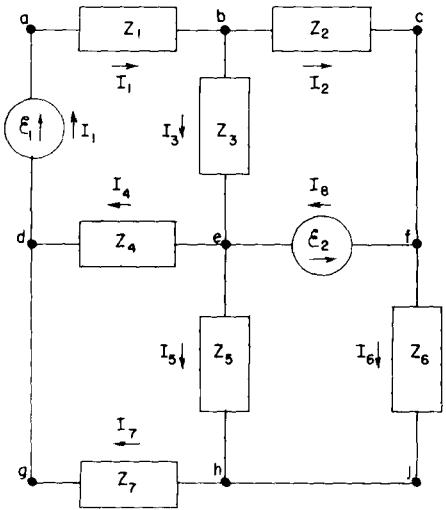


Fig. 22-11. Analyzing a circuit with Kirchhoff's rules.

Similarly, for the node labeled e we would have the current equation

$$I_3 - I_4 + I_8 - I_5 = 0.$$

For the circuit shown there are five such current equations. It turns out, however, that any one of these equations can be derived from the other four; there are, therefore, only four independent current equations. We thus have a total of eight independent, linear equations: the four voltage equations and the four current equations. With these eight equations we can solve for the eight unknown currents. Once the currents are known the circuit is solved. The voltage drop across any element is given by the current through that element times its impedance (or, in the case of the voltage sources, it is already known).

We have seen that when we write the current equations, we get one equation which is not independent of the others. Generally it is also possible to write down too many voltage equations. For example, in the circuit of Fig. 22-11, although we have considered only the four small loops, there are a large number of other loops for which we could write the voltage equation. There is, for example, the loop along the path $abcfeda$. There is another loop which follows the path $abcfhgda$. You can see that there are many loops. In analyzing complicated circuits it is very easy to get too many equations. There are rules which tell us how to proceed so that only the minimum number of equations is written down, but usually with a little thought it is possible to see how to get the right number of equations in the simplest form. Besides, writing an extra equation or two doesn't do any harm. They will not lead to any wrong answers, only perhaps a little unnecessary algebra.

In Chapter 25 of Vol. I we showed that if the two impedances z_1 and z_2 are in *series*, they are equivalent to a single impedance z_s given by

$$z_s = z_1 + z_2. \quad (22.18)$$

We also showed that if the two impedances are connected in *parallel*, they are equivalent to the single impedance z_p given by

$$z_p = \frac{1}{(1/z_1) + (1/z_2)} = \frac{z_1 z_2}{z_1 + z_2}. \quad (22.19)$$

If you look back you will see that in deriving these results we were in effect making use of Kirchhoff's rules. It is often possible to analyze a complicated circuit by repeated application of the formulas for series and parallel impedances. For instance, the circuit of Fig. 22-12 can be analyzed that way. First, the impedances z_4 and z_5 can be replaced by their parallel equivalent, and so also can z_6 and z_7 . Then the impedance z_2 can be combined with the parallel equivalent of z_6 and z_7 by the series rule. Proceeding in this way, the whole circuit can be reduced to a generator in series with a single impedance Z . The current through the generator is then just \mathcal{E}/Z . Then by working backward one can solve for the currents in each of the impedances.

There are, however, quite simple circuits which cannot be analyzed by this method, as for example the circuit of Fig. 22-13. To analyze this circuit we must

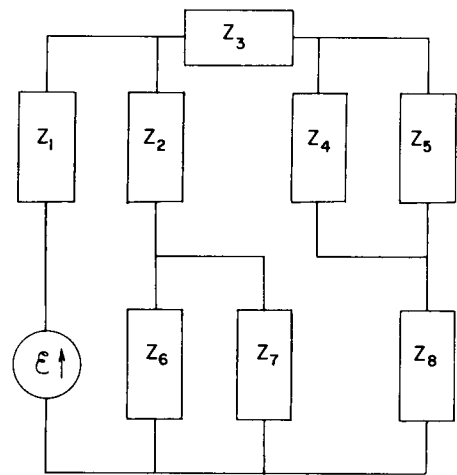


Fig. 22-12. A circuit which can be analyzed in terms of series and parallel combinations.

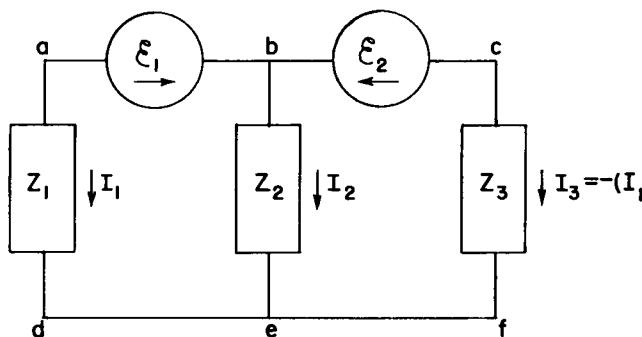


Fig. 22-13. A circuit that cannot be analyzed in terms of series and parallel combinations.

write down the current and voltage equations from Kirchhoff's rules. Let's do it. There is just one current equation:

$$I_1 + I_2 + I_3 = 0,$$

so we know immediately that

$$I_3 = -(I_1 + I_2).$$

We can save ourselves some algebra if we immediately make use of this result in writing the voltage equations. For this circuit there are two independent voltage equations; they are

$$-\varepsilon_1 + I_2 z_2 - I_1 z_1 = 0$$

and

$$\varepsilon_2 - (I_1 + I_2)z_3 - I_2 z_2 = 0.$$

There are two equations and two unknown currents. Solving these equations for I_1 and I_2 , we get

$$I_1 = \frac{z_2 \varepsilon_2 - (z_2 + z_3) \varepsilon_1}{z_1(z_2 + z_3) + z_2 z_3} \quad (22.20)$$

and

$$I_2 = \frac{z_1 \varepsilon_2 + z_3 \varepsilon_1}{z_1(z_2 + z_3) + z_2 z_3}. \quad (22.21)$$

The third current is obtained from the sum of these two.

Another example of a circuit that cannot be analyzed by using the rules for series and parallel impedance is shown in Fig. 22-14. Such a circuit is called a "bridge." It appears in many instruments used for measuring impedances. With such a circuit one is usually interested in the question: How must the various impedances be related if the current through the impedance z_3 is to be zero? We leave it for you to find the conditions for which this is so.

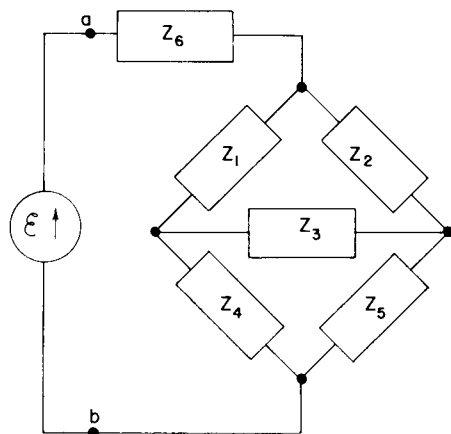


Fig. 22-14. A bridge circuit.

22-4 Equivalent circuits

Suppose we connect a generator ε to a circuit containing some complicated interconnection of impedances, as indicated schematically in Fig. 22-15(a). All of the equations we get from Kirchhoff's rules are linear, so when we solve them for the current I through the generator, we will get that I is proportional to ε . We can write

$$I = \frac{\varepsilon}{z_{\text{eff}}},$$

where now z_{eff} is some complex number, an algebraic function of all the elements in the circuit. (If the circuit contains no generators other than the one shown, there is no additional term independent of ε .) But this equation is just what we would write for the circuit of Fig. 22-15(b). So long as we are interested only in what happens *to the left* of the two terminals a and b , the two circuits of Fig. 22-15 are *equivalent*. We can, therefore, make the general statement that *any* two-terminal network of passive elements can be replaced by a single impedance z_{eff} without changing the currents and voltages in the rest of the circuit. This statement is, of course, just a remark about what comes out of Kirchhoff's rules—and ultimately from the linearity of Maxwell's equations.

The idea can be generalized to a circuit that contains generators as well as impedances. Suppose we look at such a circuit "from the point of view" of one of the impedances, which we will call z_n , as in Fig. 22-16(a). If we were to solve the equation for the whole circuit, we would find that the voltage V_n between the two terminals a and b is a linear function of I , which we can write

$$V_n = A - BI_n, \quad (22.22)$$

where A and B depend on the generators and impedances in the circuit to the left

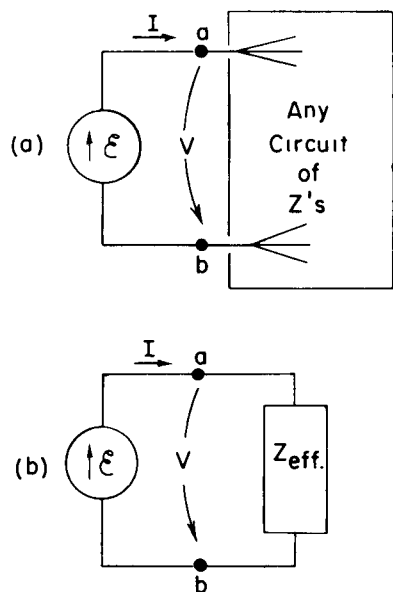


Fig. 22-15. Any two-terminal network of passive elements is equivalent to an effective impedance.

of the terminals. For instance, for the circuit of Fig. 22-13, we find $V_1 = I_1 z_1$. This can be written (by rearranging Eq. (22.20)) as

$$V_1 = \left[\left(\frac{z_2}{z_2 + z_3} \right) \mathcal{E}_2 - \mathcal{E}_1 \right] - \frac{z_2 z_3}{z_2 + z_3} I_1. \quad (22.23)$$

The complete solution is then obtained by combining this equation with the one for the impedance z_1 , namely, $V_1 = I_1 z_1$, or in the general case, by combining Eq. (22.22) with

$$V_n = I_n z_n.$$

If now we consider that z_n is attached to a simple series circuit of a generator and a current, as in Fig. 22-15(b), the equation corresponding to Eq. (22.22) is

$$V_n = \mathcal{E}_{\text{eff}} - I_n z_{\text{eff}},$$

which is identical to Eq. (22.22) provided we set $\mathcal{E}_{\text{eff}} = A$ and $z_{\text{eff}} = B$. So if we are interested only in what happens to the right of the terminals a and b , the arbitrary circuit of Fig. 22-16 can always be replaced by an equivalent combination of a generator in series with an impedance.

22-5 Energy

We have seen that to build up the current I in an inductance, the energy $U = \frac{1}{2}LI^2$ must be provided by the external circuit. When the current falls back to zero, this energy is delivered back to the external circuit. There is no energy-loss mechanism in an ideal inductance. When there is an alternating current through an inductance, energy flows back and forth between it and the rest of the circuit, but the *average* rate at which energy is delivered to the circuit is zero. We say that an inductance is a *nondissipative* element; no electrical energy is dissipated—that is, “lost”—in it.

Similarly, the energy of a condenser, $U = \frac{1}{2}CV^2$, is returned to the external circuit when a condenser is discharged. When a condenser is in an AC circuit energy flows in and out of it, but the net energy flow in each cycle is zero. An ideal condenser is also a nondissipative element.

We know that an emf is a source of energy. When a current I flows in the direction of the emf, energy is delivered to the external circuit at the rate $dU/dt = \mathcal{E}I$. If current is driven *against* the emf—by other generators in the circuit—the emf will absorb energy at the rate $\mathcal{E}I$; since I is negative, dU/dt will also be negative.

If a generator is connected to a resistor R , the current through the resistor is $I = \mathcal{E}/R$. The energy being supplied by the generator at the rate $\mathcal{E}I$ is being absorbed by the resistor. This energy goes into heat in the resistor and is lost from the electrical energy of the circuit. We say that electrical energy is *dissipated* in a resistor. The rate at which energy is dissipated in a resistor is $dU/dt = RI^2$.

In an AC circuit the average rate of energy lost to a resistor is the average of RI^2 over one cycle. Since $I = \hat{I}e^{i\omega t}$ —by which we really mean that I varies as $\cos \omega t$ —the average of I^2 over one cycle is $|\hat{I}|^2/2$, since the peak current is $|\hat{I}|$ and the average of $\cos^2 \omega t$ is $1/2$.

What about the energy loss when a generator is connected to an arbitrary impedance z ? (By “loss” we mean, of course, conversion of electrical energy into thermal energy.) Any impedance z can be written as the sum of its real and imaginary parts. That is,

$$z = R + iX, \quad (22.24)$$

where R and X are real numbers. From the point of view of equivalent circuits we can say that any impedance is equivalent to a resistance in series with a pure imaginary impedance—called a *reactance*—as shown in Fig. 22-17.

We have seen earlier that any circuit that contains only L 's and C 's has an impedance that is a pure imaginary number. Since there is no energy loss into any of the L 's and C 's on the average, a pure reactance containing only L 's and C 's will have no energy loss. We can see that this must be true in general for a reactance.

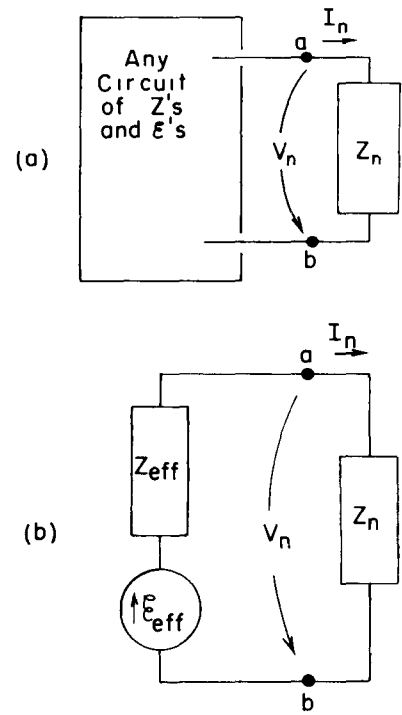


Fig. 22-16. Any two-terminal network can be replaced by a generator in series with an impedance.

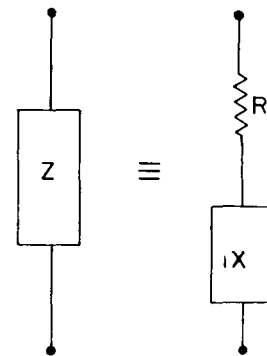


Fig. 22-17. Any impedance is equivalent to a series combination of a pure resistance and a pure reactance.

If a generator with the emf \mathcal{E} is connected to the impedance z of Fig. 22-17, the emf must be related to the current I from the generator by

$$\mathcal{E} = I(R + iX). \quad (22.25)$$

To find the average rate at which energy is delivered, we want the average of the product $\mathcal{E}I$. Now we must be careful. When dealing with such products, we must deal with the real quantities $\mathcal{E}(t)$ and $I(t)$. (The real parts of the complex functions will represent the actual physical quantities only when we have *linear* equations; now we are concerned with *products*, which are certainly not linear.)

Suppose we choose our origin of t so that the amplitude \hat{I} is a real number, let's say I_0 ; then the actual time variation I is given by

$$I = I_0 \cos \omega t.$$

The emf of Eq. (22.25) is the real part of

$$I_0 e^{i\omega t} (R + iX)$$

or

$$\mathcal{E} = I_0 R \cos \omega t - I_0 X \sin \omega t. \quad (22.26)$$

The two terms in Eq. (22.26) represent the voltage drops across R and X in Fig. 22-17. We see that the voltage drop across the resistance is *in phase* with the current, while the voltage drop across the purely reactive part is *out of phase* with the current.

The *average rate* of energy loss, $\langle P \rangle_{av}$, from the generator is the integral of the product $\mathcal{E}I$ over one cycle divided by the period T ; in other words,

$$\langle P \rangle_{av} = \frac{1}{T} \int_0^T \mathcal{E}I dt = \frac{1}{T} \int_0^T I_0^2 R \cos^2 \omega t dt - \frac{1}{T} \int_0^T I_0^2 X \cos \omega t \sin \omega t dt.$$

The first integral is $\frac{1}{2} I_0^2 R$, and the second integral is zero. So the average energy loss in an impedance $z = R + iX$ depends only on the real part of z , and is $I_0^2 R/2$, which is in agreement with our earlier result for the energy loss in a resistor. There is no energy loss in the reactive part.

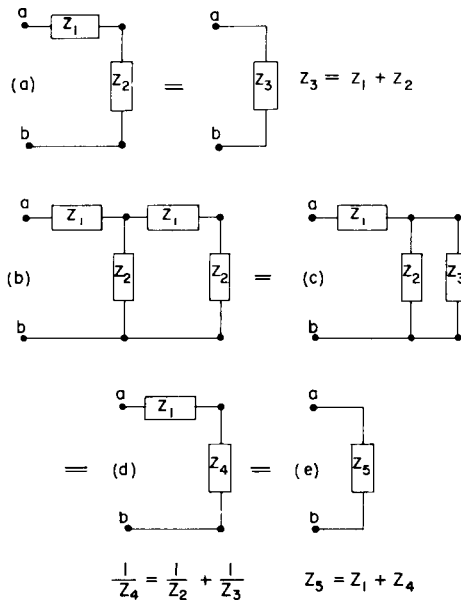


Fig. 22-18. The effective impedance of a ladder.

22-6 A ladder network

We would like now to consider an interesting circuit which can be analyzed in terms of series and parallel combinations. Suppose we start with the circuit of Fig. 22-18(a). We can see right away that the impedance from terminal a to terminal b is simply $z_1 + z_2$. Now let's take a little harder circuit, the one shown in Fig. 22-18(b). We could analyze this circuit using Kirchhoff's rules, but it is also easy to handle with series and parallel combinations. We can replace the two impedances on the right-hand end by a single impedance $z_3 = z_1 + z_2$, as in part (c) of the figure. Then the two impedances z_2 and z_3 can be replaced by their equivalent parallel impedance z_4 , as shown in part (d) of the figure. Finally, z_1 and z_4 are equivalent to a single impedance z_5 , as shown in part (e).

Now we may ask an amusing question: What would happen in the network of Fig. 22-18(b) we kept on adding more sections *forever*—as we indicate by the dashed lines in Fig. 22-19(a)? Can we solve such an infinite network? Well, that's

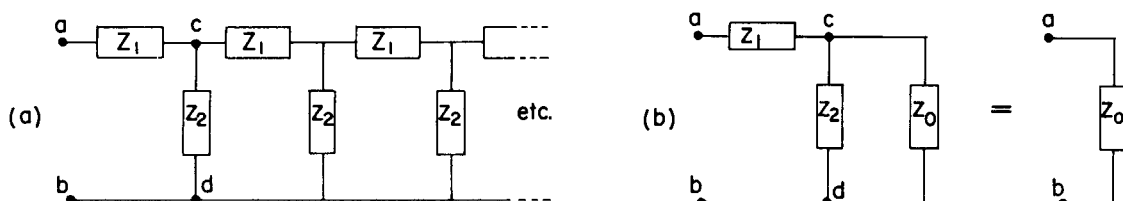


Fig. 22-19. The effective impedance of an infinite ladder.

not so hard. First, we notice that such an infinite network is unchanged if we add one more section at the “front” end. Surely, if we add one more section to an infinite network it is still the same infinite network. Suppose we call the impedance between the two terminals a and b of the infinite network z_0 ; then the impedance of all the stuff to the right of the two terminals c and d is also z_0 . Therefore, so far as the front end is concerned, we can represent the network as shown in Fig. 22-19(b). Combining the parallel combinations $z_2 z_0$ and adding the result in series with z_1 , we can immediately write down the impedance of this combination:

$$z = z_1 + \frac{1}{(1/z_2) + (1/z_0)} \quad \text{or} \quad z = z_1 + \frac{z_2 z_0}{z_2 + z_0}.$$

But this impedance is also equal to z_0 , so we have the equation

$$z_0 = z_1 + \frac{z_2 z_0}{z_2 + z_0}.$$

We can solve for z_0 to get

$$z_0 = \frac{z_1}{2} + \sqrt{(z_1^2/4) + z_1 z_2}. \quad (22.27)$$

So we have found the solution for the impedance of an infinite ladder of repeated series and parallel impedances. The impedance z_0 is called the *characteristic impedance* of such an infinite network.

Let's now consider a specific example in which the series element is an inductance L and the shunt element is a capacitance C , as shown in Fig. 22-20(a). In this case we find the impedance of the infinite network by setting $z_1 = i\omega L$ and $z_2 = 1/i\omega C$. Notice that the first term, $z_1/2$, in Eq. (22.27) is just one-half the impedance of the first element. It would therefore seem more natural, or at least somewhat simpler, if we were to draw our infinite network as shown in Fig. 22-20(b). Looking at the infinite network from the terminal a' we would see the characteristic impedance

$$z_0 = \sqrt{(L/C) - (\omega^2 L^2/4)}. \quad (22.28)$$

Now there are two interesting cases, depending on the frequency ω . If ω^2 is less than $4/LC$, the second term in the radical will be smaller than the first, and the impedance z_0 will be a real number. On the other hand, if ω^2 is greater than $4/LC$ the impedance z_0 will be a pure imaginary number which we can write as

$$z_0 = i\sqrt{(\omega^2 L^2/4) - (L/C)}.$$

We have said earlier that a circuit which contains only imaginary impedances, such as inductances and capacitances, will have an impedance which is purely imaginary. How can it be then that for the circuit we are now studying—which has only L 's and C 's—the impedance is a pure resistance for frequencies below $\sqrt{4/LC}$? For higher frequencies the impedance is purely imaginary, in agreement with our earlier statement. For lower frequencies the impedance is a pure resistance and will therefore absorb energy. But how can the circuit continuously absorb energy, as a resistance does, if it is made only of inductances and capacitances? *Answer:* Because there is an infinite number of inductances and capacitances, so that when a source is connected to the circuit, it supplies energy to the first inductance and capacitance, then to the second, to the third, and so on. In a circuit of this kind, energy is continually absorbed from the generator at a constant rate and flows constantly out into the network, supplying energy which is stored in the inductances and capacitances down the line.

This idea suggests an interesting point about what is happening in the circuit. We would expect that if we connect a source to the front end, the effects of this source will be propagated through the network toward the infinite end. The propagation of the waves down the line is much like the radiation from an antenna which absorbs energy from its driving source; that is, we expect such a propagation to occur when the impedance is real, which occurs if ω is less than $\sqrt{4/LC}$. But when the impedance is purely imaginary, which happens for ω greater than $\sqrt{4/LC}$, we would not expect to see any such propagation.

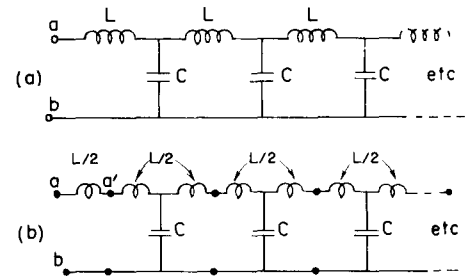


Fig. 22-20. An L - C ladder drawn in two equivalent ways.

22-7 Filters

We saw in the last section that the infinite ladder network of Fig. 22-20 absorbs energy continuously if it is driven at a frequency below a certain critical frequency $\sqrt{4/LC}$, which we will call the *cutoff frequency* ω_0 . We suggested that this effect could be understood in terms of a continuous transport of energy down the line. On the other hand, at high frequencies, for $\omega > \omega_0$, there is no continuous absorption of energy; we should then expect that perhaps the currents don't "penetrate" very far down the line. Let's see whether these ideas are right.

Suppose we have the front end of the ladder connected to some AC generator and we ask what the voltage looks like at, say, the 754th section of the ladder. Since the network is infinite, whatever happens to the voltage from one section to the next is always the same; so let's just look at what happens when we go from some section, say the n th to the next. We will define the currents I_n and voltages V_n as shown in Fig. 22-21(a).

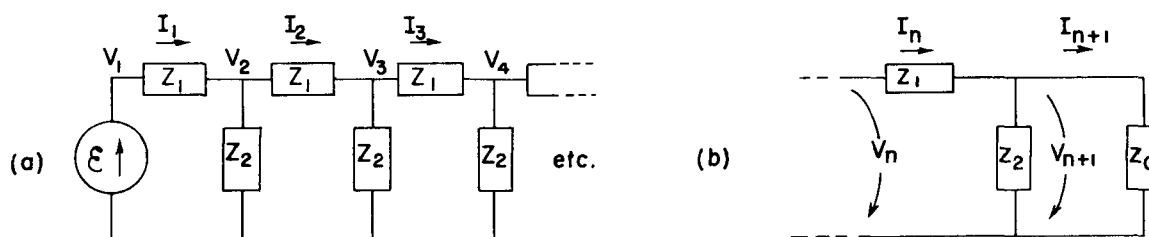


Fig. 22-21. Finding the propagation factor of a ladder.

We can get the voltage V_{n+1} from V_n by remembering that we can always replace the rest of the ladder after the n th section by its characteristic impedance z_0 ; then we need only analyze the circuit of Fig. 22-21(b). First, we notice that any V_n , since it is across z_0 , must equal $I_n z_0$. Also, the difference between V_n and V_{n+1} is just $I_n z_1$:

$$V_n - V_{n+1} = I_n z_1 = V_n \frac{z_1}{z_0}$$

So we get the ratio

$$\frac{V_{n+1}}{V_n} = 1 - \frac{z_1}{z_0} = \frac{z_0 - z_1}{z_0}.$$

We can call this ratio the *propagation factor* for one section of the ladder; we'll call it α . It is, of course, the same for all sections:

$$\alpha = \frac{z_0 - z_1}{z_0}. \quad (22.29)$$

The voltage after the n th section is then

$$V_n = \alpha^n \epsilon. \quad (22.30)$$

You can now find the voltage after 754 sections; it is just α to the 754th power times ϵ .

Suppose we see what α is like for the L - C ladder of Fig. 22-20(a). Using z_0 from Eq. (22.27), and $z_1 = i\omega L$, we get

$$\alpha = \frac{\sqrt{(L/C) - (\omega^2 L^2/4)} - i(\omega L/2)}{\sqrt{(L/C) - (\omega^2 L^2/4)} + i(\omega L/2)} \quad (22.31)$$

If the driving frequency is below the cutoff frequency $\omega_0 = \sqrt{4/LC}$, the radical is a real number, and the magnitudes of the complex numbers in the numerator and denominator are equal. Therefore, the magnitude of α is one; we can write

$$\alpha = e^{i\delta},$$

which means that the magnitude of the voltage is the same at every section, only

its phase changes. The phase change δ is, in fact, a negative number and represents the “delay” of the voltage as it passes along the network.

For frequencies above the cutoff frequency ω_0 it is better to factor out an i from the numerator and denominator of Eq. (22-31) and rewrite it as

$$\alpha = \frac{\sqrt{(\omega^2 L^2/4) - (L/C)} - (\omega L/2)}{\sqrt{(\omega^2 L^2/4) - (L/C)} + (\omega L/2)}. \quad (22-32)$$

The propagation factor α is now a *real* number, and a number *less than one*. That means that the voltage at any section is always less than the voltage at the preceding section by the factor α . For any frequency above ω_0 , the voltage dies away rapidly as we go along the network. A plot of the absolute value of α as a function of frequency looks like the graph in Fig. 22-22.

We see that the behavior of α , both above and below ω_0 , agrees with our interpretation that the network propagates energy for $\omega < \omega_0$ and blocks it for $\omega > \omega_0$. We say that the network “passes” low frequencies and “rejects” or “filters out” the high frequencies. Any network designed to have its characteristics vary in a prescribed way with frequency is called a “filter.” We have been analyzing a “low-pass filter.”

You may be wondering why all this discussion of an infinite network which obviously cannot actually occur. The point is that the same characteristics are found in a finite network if we finish it off at the end with an impedance equal to the characteristic impedance z_0 . Now in practice it is not possible to *exactly* reproduce the characteristic impedance with a few simple elements—like R ’s, L ’s, and C ’s. But it is often possible to do so with a fair approximation for a certain range of frequencies. In this way one can make a finite filter network whose properties are very nearly the same as those for the infinite case. For instance, the L - C ladder behaves much as we have described it if it is terminated in the pure resistance $R = \sqrt{L/C}$.

If in our L - C ladder we interchange the positions of the L ’s and C ’s, to make the ladder shown in Fig. 22-23(a), we can have a filter that propagates *high* frequencies and rejects *low* frequencies. It is easy to see what happens with this network by using the results we already have. You will notice that whenever we change an L to a C and *vice versa*, we also change every $i\omega$ to $1/i\omega$. So whatever happened at ω before will now happen at $1/\omega$. In particular, we can see how α will vary with frequency by using Fig. 22-22 and changing the label on the axis to $1/\omega$, as we have done in Fig. 22-23(b).

The low-pass and high-pass filters we have described have various technical applications. An L - C low-pass filter is often used as a “smoothing” filter in a DC power supply. If we want to manufacture DC power from an AC source, we begin with a rectifier which permits current to flow only in one direction. From the rectifier we get a series of pulses that look like the function $V(t)$ shown in Fig. 22-24, which is lousy DC, because it wobbles up and down. Suppose we would like a nice pure DC, such as a battery provides. We can come close to that by putting a low-pass filter between the rectifier and the load.

We know from Chapter 50 of Vol. I that the time function in Fig. 22-24 can be represented as a superposition of a constant voltage plus a sine wave, plus a higher-frequency sine wave, plus a still higher-frequency sine wave, etc.—by a Fourier series. If our filter is linear (if, as we have been assuming, the L ’s and C ’s don’t vary with the currents or voltages) then what comes out of the filter is the superposition of the outputs for each component at the input. If we arrange that the cutoff frequency ω_0 of our filter is well below the lowest frequency in the function $V(t)$, the DC (for which $\omega = 0$) goes through fine, but the amplitude of the first harmonic will be cut down a lot. And amplitudes of the higher harmonics will be cut down even more. So we can get the output as smooth as we wish, depending only on how many filter sections we are willing to buy.

A high-pass filter is used if one wants to reject certain low frequencies. For instance, in a phonograph amplifier a high-pass filter may be used to let the music

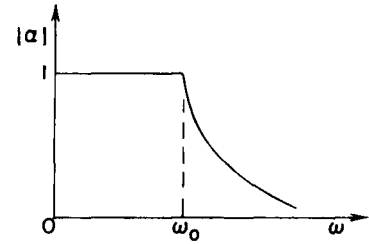
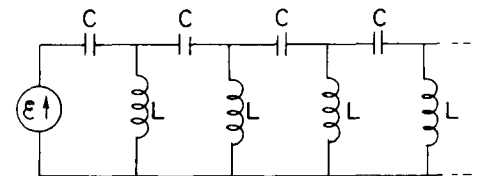
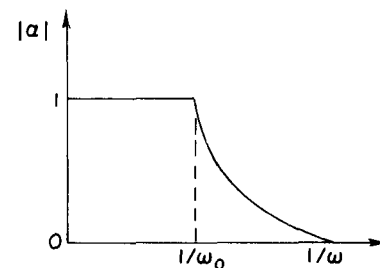


Fig. 22-22. The propagation factor of a section of an L - C ladder



(a)



(b)

Fig. 22-23. (a) A high-pass filter; (b) its propagation factor as a function of $1/\omega$.

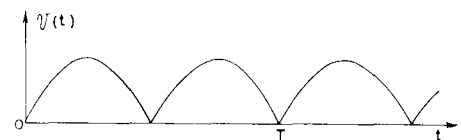


Fig. 22-24. The output voltage of a full-wave rectifier.

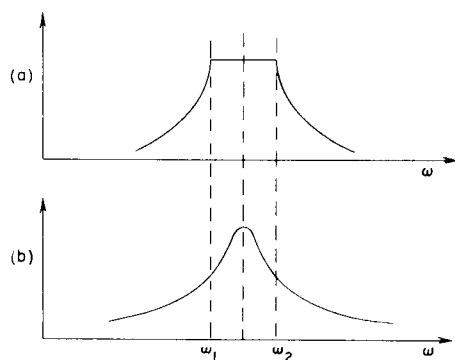


Fig. 22-25. (a) A band-pass filter. (b) A simple resonant filter.

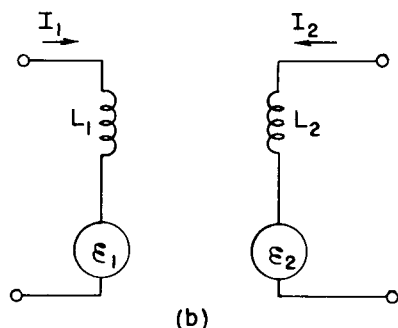
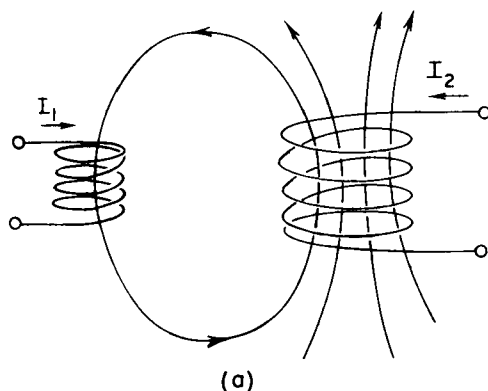


Fig. 22-26. Equivalent circuit of a mutual inductance.

through, while keeping out the low-pitched rumbling from the motor of the turntable.

It is also possible to make “band-pass” filters that reject frequencies below some frequency ω_1 and above another frequency ω_2 (greater than ω_1), but pass the frequencies between ω_1 and ω_2 . This can be done simply by putting together a high-pass and a low-pass filter, but it is more usually done by making a ladder in which the impedances z_1 and z_2 are more complicated—being each a combination of L ’s and C ’s. Such a band-pass filter might have a propagation constant like that shown in Fig. 22-25(a). It might be used, for example, in separating signals that occupy only an interval of frequencies, such as each of the many voice channels in a high-frequency telephone cable, or the modulated carrier of a radio transmission.

We have seen in Chapter 25 of Vol. I that such filtering can also be done using the selectivity of an ordinary resonance curve, which we have drawn for comparison in Fig. 22-25(b). But the resonant filter is not as good for some purposes as the band-pass filter. You will remember (Chapter 48, Vol. I) that when a carrier of frequency ω_c is modulated with a “signal” frequency ω_s , the total signal contains not only the carrier frequency but also the two side-band frequencies $\omega_c + \omega_s$ and $\omega_c - \omega_s$. With a resonant filter, these side-bands are always attenuated somewhat, and the attenuation is more, the higher the signal frequency, as you can see from the figure. So there is a poor “frequency response.” The higher musical tones don’t get through. But if the filtering is done with a band-pass filter designed so that the width $\omega_2 - \omega_1$ is at least twice the highest signal frequency, the frequency response will be “flat” for the signals wanted.

We want to make one more point about the ladder filter: the L - C ladder of Fig. 22-20 is also an approximate representation of a transmission line. If we have a long conductor that runs parallel to another conductor—such as a wire in a coaxial cable, or a wire suspended above the earth—there will be some capacitance between the two conductors and also some inductance due to the magnetic field between them. If we imagine the line as broken up into small lengths Δl , each length will look like one section of the L - C ladder with a series inductance ΔL and a shunt capacitance ΔC . We can then use our results for the ladder filter. If we take the limit as Δl goes to zero, we have a good description of the transmission line. Notice that as Δl is made smaller and smaller, both ΔL and ΔC decrease, but in the same proportion, so that the ratio $\Delta L/\Delta C$ remains constant. So if we take the limit of Eq. (22.28) as ΔL and ΔC go to zero, we find that the characteristic impedance z_0 is a pure resistance whose magnitude is $\sqrt{\Delta L/\Delta C}$. We can also write the ratio $\Delta L/\Delta C$ as L_0/C_0 , where L_0 and C_0 are the inductance and capacitance of a unit length of the line; then we have

$$z_0 = \sqrt{\frac{L_0}{C_0}}. \quad (22.33)$$

You will also notice that as ΔL and ΔC go to zero, the cutoff frequency $\omega_0 = \sqrt{4/LC}$ goes to infinity. There is no cutoff frequency for an ideal transmission line.

22-8 Other circuit elements

We have so far defined only the ideal circuit impedances—the inductance, the capacitance, and the resistance—as well as the ideal voltage generator. We want now to show that other elements, such as mutual inductances or transistors or vacuum tubes, can be described by using only the same basic elements. Suppose that we have two coils and that on purpose, or otherwise, some flux from one of the coils links the other, as shown in Fig. 22-26(a). Then the two coils will have a mutual inductance M such that when the current varies in one of the coils, there will be a voltage generated in the other. Can we take into account such an effect in our equivalent circuits? We can in the following way. We have seen that the

induced emf's in each of two interacting coils can be written as the sum of two parts:

$$\begin{aligned}\varepsilon_1 &= -L_1 \frac{dI_1}{dt} \pm M \frac{dI_2}{dt}, \\ \varepsilon_2 &= -L_2 \frac{dI_2}{dt} \pm M \frac{dI_1}{dt}.\end{aligned}\quad (22.34)$$

The first term comes from the self-inductance of the coil, and the second term comes from its mutual inductance with the other coil. The sign of the second term can be plus or minus, depending on the way the flux from one coil links the other. Making the same approximations we used in describing an ideal inductance, we would say that the potential difference across the terminals of each coil is equal to the electromotive force in the coil. Then the two equations of (22.34) are the same as the ones we would get from the circuit of Fig. 22-26(b), provided the electromotive force in each of the two circuits shown depends on the current in the opposite circuit according to the relations

$$\varepsilon_1 = \pm i\omega MI_2, \quad \varepsilon_2 = \pm i\omega MI_1. \quad (22.35)$$

So what we can do is represent the effect of the self-inductance in a normal way but replace the effect of the mutual inductance by an auxiliary ideal voltage generator. We must in addition, of course, have the equation that relates this emf to the current in some other part of the circuit; but so long as this equation is linear, we have just added more linear equations to our circuit equations, and all of our earlier conclusions about equivalent circuits and so forth are still correct.

In addition to mutual inductances there may also be mutual capacitances. So far, when we have talked about condensers we have always imagined that there were only two electrodes, but in many situations, for example in a vacuum tube, there may be many electrodes close to each other. If we put an electric charge on any one of the electrodes, its electric field will induce charges on each of the other electrodes and affect its potential. As an example, consider the arrangement of four plates shown in Fig. 22-27(a). Suppose these four plates are connected to external circuits by means of the wires *A*, *B*, *C*, and *D*. So long as we are only worried about electrostatic effects, the equivalent circuit of such an arrangement of electrodes is as shown in part (b) of the figure. The electrostatic interaction of any electrode with each of the others is equivalent to a capacity between the two electrodes.

Finally, let's consider how we should represent such complicated devices as transistors and radio tubes in an AC circuit. We should point out at the start that such devices are often operated in such a way that the relationship between the currents and voltages is not at all linear. In such cases, those statements we have made which depend on the linearity of equations are, of course, no longer correct. On the other hand, in many applications the operating characteristics are sufficiently linear that we may consider the transistors and tubes to be linear devices. By this we mean that the alternating currents in, say, the plate of a vacuum tube are linearly proportional to the voltages that appear on the other electrodes, say the grid voltage and the plate voltage. When we have such linear relationships, we can incorporate the device into our equivalent circuit representation.

As in the case of the mutual inductance, our representation will have to include auxiliary voltage generators which describe the influence of the voltages or currents in one part of the device on the currents or voltages in another part. For example, the plate circuit of a triode can usually be represented by a resistance in series with an ideal voltage generator whose source strength is proportional to the grid voltage. We get the equivalent circuit shown in Fig. 22-28.* Similarly, the collector circuit

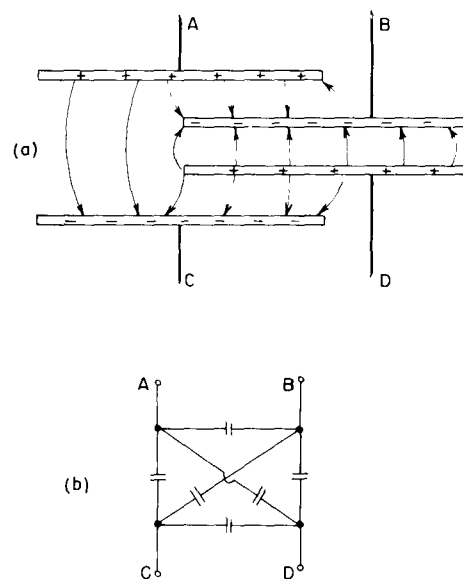


Fig. 22-27. Equivalent circuit of mutual capacitance.

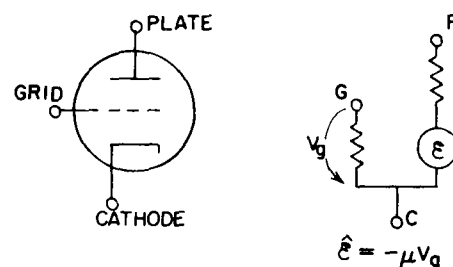


Fig. 22-28. A low-frequency equivalent circuit of a vacuum triode.

* The equivalent circuit shown is correct only for low frequencies. For high frequencies the equivalent circuit gets much more complicated and will include various so-called "parasitic" capacitances and inductances.

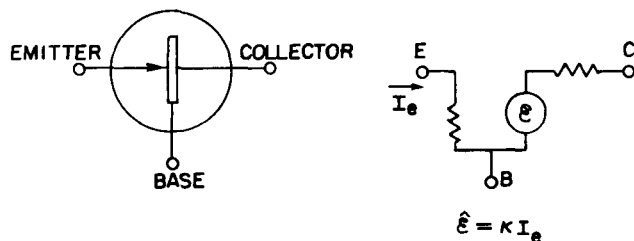


Fig. 22-29. A low-frequency equivalent circuit of a transistor.

of a transistor is conveniently represented as a resistor in series with an ideal voltage generator whose source strength is proportional to the current from the emitter to the base of the transistor. The equivalent circuit is then like that in Fig. 22-29. So long as the equations which describe the operation are linear, we can use such representations for tubes or transistors. Then, when they are incorporated in a complicated network, our general conclusions about the equivalent representation of any arbitrary connection of elements is still valid.

There is one remarkable thing about transistor and radio tube circuits which is different from circuits containing only impedances: the real part of the effective impedance z_{eff} can become negative. We have seen that the real part of z represents the loss of energy. But it is the important characteristic of transistors and tubes that they *supply* energy to the circuit. (Of course they don't just "make" energy; they take energy from the DC circuits of the power supplies and convert it into AC energy.) So it is possible to have a circuit with a negative resistance. Such a circuit has the property that if you connect it to an impedance with a positive real part, i.e., a positive resistance, and arrange matters so that the sum of the two real parts is exactly zero, then there is no dissipation in the combined circuit. If there is no loss of energy, any alternating voltage once started will remain forever. This is the basic idea behind the operation of an oscillator or signal generator which can be used as a source of alternating voltage at any desired frequency.

Cavity Resonators

23-1 Real circuit elements

When looked at from any one pair of terminals, any arbitrary circuit made up of ideal impedances and generators is, at any given frequency, equivalent to a generator \mathcal{E} in series with an impedance z . That comes about because if we put a voltage V across the terminals and solve all the equations to find the current I , we must get a linear relation between the current and the voltage. Since all the equations are linear, the result for I must also depend only linearly on V . The most general linear form can be expressed as

$$I = \frac{1}{z} (V - \mathcal{E}). \quad (23.1)$$

In general, both z and \mathcal{E} may depend in some complicated way on the frequency ω . Equation (23.1), however, is the relation we would get if behind the two terminals there was just the generator $\mathcal{E}(\omega)$ in series with the impedance $z(\omega)$.

There is also the opposite kind of question. If we have any electromagnetic device at all with two terminals and we *measure* the relation between I and V to determine \mathcal{E} and z as functions of frequency, can we find a combination of our ideal elements that is equivalent to the internal impedance z ? The answer is that for any reasonable—that is, physically meaningful—function $z(\omega)$, it is possible to *approximate* the situation to as high an accuracy as you wish with a circuit containing a finite set of ideal elements. We don't want to consider the general problem now but only look at what might be expected from physical arguments for a few cases.

If we think of a real resistor, we know that the current through it will produce a magnetic field. So any real resistor should also have some inductance. Also, when a resistor has a potential difference across it, there must be charges on the ends of the resistor to produce the necessary electric fields. As the voltage changes, the charges will change in proportion, so the resistor will also have some capacitance. We expect that a *real* resistor might have the equivalent circuit shown in Fig. 23-1. In a well-designed resistor, the so-called “parasitic” elements L and C are small, so that at the frequencies for which it is intended, ωL is much less than R , and $1/\omega C$ is much greater than R . It may therefore be possible to neglect them. As the frequency is raised, however, they will eventually become important, and a resistor begins to look like a resonant circuit.

A real inductance is also not equal to the idealized inductance, whose impedance is $i\omega L$. A real coil of wire will have some resistance, so at low frequencies the coil is really equivalent to an inductance in series with some resistance, as shown in Fig. 23-2(a). But, you are thinking, the resistance and inductance are *together* in a real coil—the resistance is spread all along the wire, so it is mixed in with the inductance. We should probably use a circuit more like the one in Fig. 23-2(b), which has several little R 's and L 's in series. But the total impedance of such a circuit is just $\sum R + \sum i\omega L$, which is equivalent to the simpler diagram of part (a).

As we go up in frequency with a real coil, the approximation of an inductance plus a resistance is no longer very good. The charges that must build up on the wires to make the voltages will become important. It is as if there were little condensers across the turns of the coil, as sketched in Fig. 23-3(a). We might try to approximate the real coil by the circuit in Fig. 23-3(b). At low frequencies, this circuit can be imitated fairly well by the simpler one in part (c) of the figure (which is again the same resonant circuit we found for the high-frequency model of a resistor). For higher frequencies, however, the more complicated circuit of

23-1 Real circuit elements

23-2 A capacitor at high frequencies

23-3 A resonant cavity

23-4 Cavity modes

23-5 Cavities and resonant circuits

Review: Chapter 23, Vol. I, *Resonance*
Chapter 49, Vol. I, *Modes*

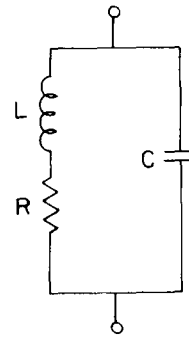


Fig. 23-1. Equivalent circuit of a real resistor.

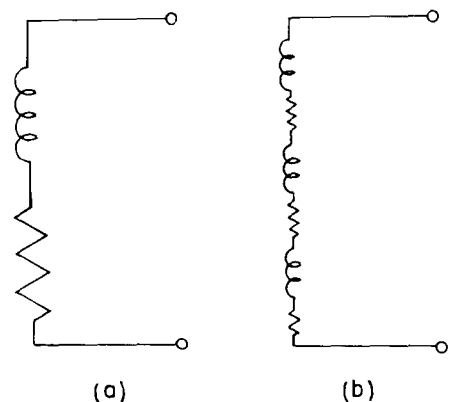


Fig. 23-2. The equivalent circuit of a real inductance at low frequencies.

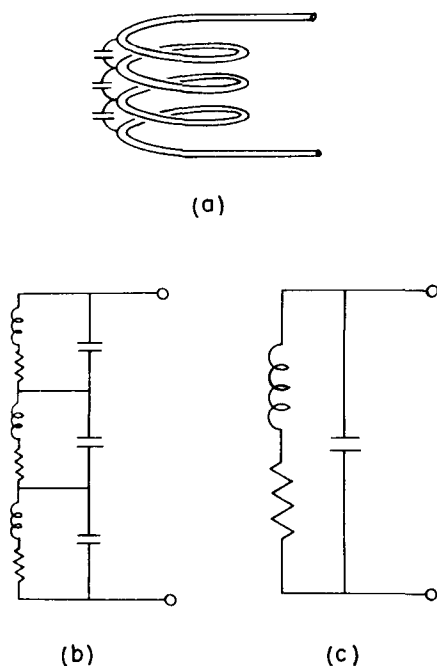


Fig. 23-3. The equivalent circuit of a real inductance at higher frequencies.

Fig. 23-3(b) is better. In fact, the more accurately you wish to represent the actual impedance of a real, physical inductance, the more ideal elements you will have to use in the artificial model of it.

Let's look a little more closely at what goes on in a real coil. The impedance of an inductance goes as ωL , so it becomes zero at low frequencies—it is a "short circuit": all we see is the resistance of the wire. As we go up in frequency, ωL soon becomes much larger than R , and the coil looks pretty much like an ideal inductance. As we go still higher, however, the capacities become important. Their impedance is proportional to $1/\omega C$, which is large for small ω . For small enough frequencies a condenser is an "open circuit," and when it is in parallel with something else, it draws no current. But at high frequencies, the current prefers to flow into the capacitance between the turns, rather than through the inductance. So the current in the coil jumps from one turn to the other and doesn't bother to go around and around where it has to buck the emf. So although we may have *intended* that the current should go around the loop, it will take the easier path—the path of least impedance.

If the subject had been one of popular interest, this effect would have been called "the high-frequency barrier," or some such name. The same kind of thing happens in all subjects. In aerodynamics, if you try to make things go faster than the speed of sound when they were designed for lower speeds, they don't work. It doesn't mean that there is a great "barrier" there, it just means that the object should be redesigned. So this coil which we designed as an "inductance" is not going to work as a good inductance, but as some other kind of thing at very high frequencies. For high frequencies, we have to find a new design.

23-2 A capacitor at high frequencies

Now we want to discuss in detail the behavior of a capacitor—a geometrically ideal capacitor—as the frequency gets larger and larger, so we can see the transition of its properties. (We prefer to use a capacitor instead of an inductance, because the geometry of a pair of plates is much less complicated than the geometry of a coil.) We consider the capacitor shown in Fig. 23-4(a), which consists of two parallel circular plates connected to an external generator by a pair of wires. If we charge the capacitor with DC, there will be a positive charge on one plate and a negative charge on the other; and there will be a uniform electric field between the plates.

Now suppose that instead of DC, we put an AC of low frequency on the plates. (We will find out later what is "low" and what is "high".) Say we connect the capacitor to a lower-frequency generator. As the voltage alternates, the positive charge on the top plate is taken off and negative charge is put on. While that is happening, the electric field disappears and then builds up in the opposite direction.

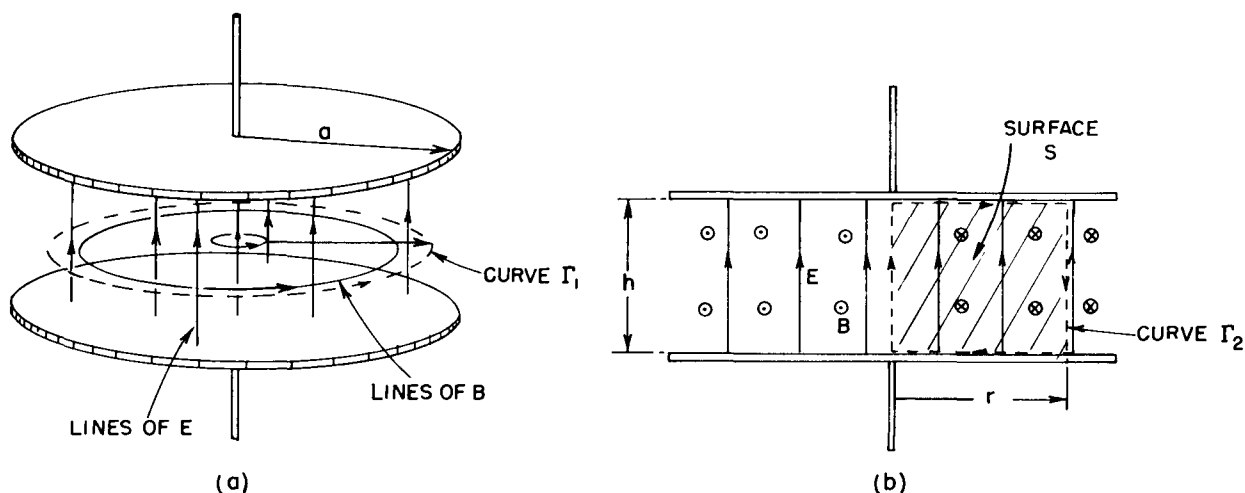


Fig. 23-4. The electric and magnetic fields between the plates of a capacitor.

As the charge sloshes back and forth slowly, the electric field follows. At each instant the electric field is uniform, as shown in Fig. 23-4(b), except for some edge effects which we are going to disregard. We can write the magnitude of the electric field as

$$E = E_0 e^{i\omega t}, \quad (23.2)$$

where E_0 is a constant.

Now will that continue to be right as the frequency goes up? No, because as the electric field is going up and down, there is a flux of electric field through any loop like Γ_1 in Fig. 23-4(a). And, as you know, a changing electric field acts to produce a magnetic field. One of Maxwell's equations says that when there is a varying electric field, as there is here, there has got to be a line integral of the magnetic field. The integral of the magnetic field around a closed ring, multiplied by c^2 , is equal to the time rate-of-change of the electric flux through the area inside the ring (if there are no currents):

$$c^2 \oint_{\Gamma} \mathbf{B} \cdot d\mathbf{s} = \frac{\partial}{\partial t} \int_{\text{inside } \Gamma} \mathbf{E} \cdot \mathbf{n} da. \quad (23.3)$$

So how much magnetic field is there? That's not very hard. Suppose that we take the loop Γ_1 , which is a circle of radius r . We can see from symmetry that the magnetic field goes around as shown in the figure. Then the line integral of \mathbf{B} is $2\pi r B$. And, since the electric field is uniform, the flux of the electric field is simply E multiplied by πr^2 , the area of the circle:

$$c^2 B \cdot 2\pi r = \frac{\partial}{\partial t} E \cdot \pi r^2. \quad (23.4)$$

The derivative of E with respect to time is, for our alternating field, simply $i\omega E_0 e^{i\omega t}$. So we find that our capacitor has the magnetic field

$$B = \frac{i\omega r}{2c^2} E_0 e^{i\omega t}. \quad (23.5)$$

In other words, the magnetic field also oscillates and has a strength proportional to r .

What is the effect of that? When there is a magnetic field that is varying, there will be induced electric fields and the capacitor will begin to act a little bit like an inductance. As the frequency goes up, the magnetic field gets stronger; it is proportional to the rate of change of E , and so to ω . The impedance of the capacitor will no longer be simply $1/i\omega C$.

Let's continue to raise the frequency and to analyze what happens more carefully. We have a magnetic field that goes sloshing back and forth. But then the electric field cannot be uniform, as we have assumed! When there is a varying magnetic field, there must be a line integral of the electric field—because of Faraday's law. So if there is an appreciable magnetic field, as begins to happen at high frequencies, the electric field cannot be the same at all distances from the center. The electric field must change with r so that the line integral of the electric field can equal the changing flux of the magnetic field.

Let's see if we can figure out the correct electric field. We can do that by computing a "correction" to the uniform field we originally assumed for low frequencies. Let's call the uniform field E_1 , which will still be $E_0 e^{i\omega t}$, and write the correct field as

$$E = E_1 + E_2,$$

where E_2 is the correction due to the changing magnetic field. For any ω we will write the field at the center of the condenser as $E_0 e^{i\omega t}$ (thereby defining E_0), so that we have no correction at the center; $E_2 = 0$ at $r = 0$.

To find E_2 we can use the integral form of Faraday's law:

$$\oint_{\Gamma} \mathbf{E} \cdot d\mathbf{s} = - \frac{\partial}{\partial t} (\text{flux of } B).$$

The integrals are simple if we take them for the curve Γ_2 , shown in Fig. 23-4(b), which goes up along the axis, out radially the distance r along the top plate, down vertically to the bottom plate, and back to the axis. The line integral of E_1 around this curve is, of course, zero; so only E_2 contributes, and its integral is just $-E_2(r) \cdot h$, where h is the spacing between the plates. (We call E positive if it points upward.) This is equal to the rate of change of the flux of B , which we have to get by an integral over the shaded area S inside Γ_2 in Fig. 23-4(b). The flux through a vertical strip of width dr is $B(r)h \, dr$, so the total flux is

$$h \int B(r) \, dr.$$

Setting $-\partial/\partial t$ of the flux equal to the line integral of E_2 , we have

$$E_2(r) = \frac{\partial}{\partial t} \int B(r) \, dr. \quad (23.6)$$

Notice that the h cancels out, the fields don't depend on the separation of the plates.

Using Eq. (23.5) for $B(r)$, we have

$$E_2(r) = \frac{\partial}{\partial t} \frac{i\omega r^2}{4c^2} E_0 e^{i\omega t}.$$

The time derivative just brings down another factor $i\omega$; we get

$$E_2(r) = -\frac{\omega^2 r^2}{4c^2} E_0 e^{i\omega t}. \quad (23.7)$$

As we expect, the induced field tends to *reduce* the electric field farther out. The corrected field $E = E_1 + E_2$ is then

$$E = E_1 + E_2 = \left(1 - \frac{1}{4} \frac{\omega^2 r^2}{c^2}\right) E_0 e^{i\omega t}. \quad (23.8)$$

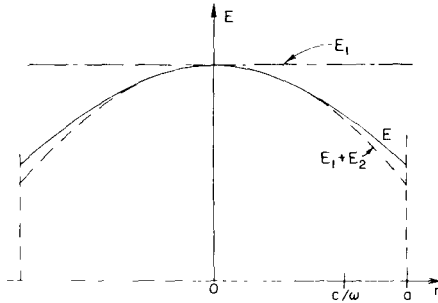


Fig. 23-5. The electric field between the capacitor plates at high frequency. (Edge effects are neglected.)

The electric field in the capacitor is no longer uniform; it has the parabolic shape shown by the broken line in Fig. 23-5. You see that our simple capacitor is getting slightly complicated.

We could now use our results to calculate the impedance of the capacitor at high frequencies. Knowing the electric field, we could compute the charges on the plates and find out how the current through the capacitor depends on the frequency ω , but we are not interested in that problem for the moment. We are more interested in seeing what happens as we continue to go up with the frequency—to see what happens at even higher frequencies. Aren't we already finished? No, because we have corrected the electric field, which means that the magnetic field we have calculated is no longer right. The magnetic field of Eq. (23.5) is approximately right, but it is only a first approximation. So let's call it B_1 . We should then rewrite Eq. (23.5) as

$$B_1 = \frac{i\omega r}{2c^2} E_0 e^{i\omega t}. \quad (23.9)$$

You will remember that this field was produced by the variation of E_1 . Now the correct magnetic field will be that produced by the total electric field $E_1 + E_2$. If we write the magnetic field as $B = B_1 + B_2$, the second term is just the additional field produced by E_2 . To find B_2 we can go through the same arguments we have used to find B_1 , the line integral of B_2 around the curve Γ_1 is equal to the rate of change of the flux of E_2 through Γ_1 . We will just have Eq. (23.4) again with B replaced by B_2 and E replaced by E_2 :

$$c^2 B_2 \cdot 2\pi r = \frac{\partial}{\partial t} (\text{flux of } E_2 \text{ through } \Gamma_1).$$

Since E_2 varies with radius, to obtain its flux we must integrate over the circular

surface inside Γ_1 . Using $2\pi r dr$ as the element of area, this integral is

$$\int_0^r E_2(r) \cdot 2\pi r dr.$$

So we get for $B_2(r)$

$$B_2(r) = \frac{1}{rc^2} \frac{\partial}{\partial t} \int E_2(r)r dr. \quad (23.10)$$

Using $E_2(r)$ from Eq. (23.7), we need the integral of $r^3 dr$, which is, of course, $r^4/4$. Our correction to the magnetic field becomes

$$B_2(r) = -\frac{i\omega^3 r^3}{16c^4} E_0 e^{i\omega t}. \quad (23.11)$$

But we are still not finished! If the magnetic field B is not the same as we first thought, then we have incorrectly computed E_2 . We must make a further correction to E , which comes from the extra magnetic field B_2 . Let's call this additional correction to the electric field E_3 . It is related to the magnetic field B_2 in the same way that E_2 was related to B_1 . We can use Eq. (23.6) all over again just by changing the subscripts:

$$E_3(r) = \frac{\partial}{\partial t} \int B_2(r) dr. \quad (23.12)$$

Using our result, Eq. (23.11), for B_2 , the new correction to the electric field is

$$E_3(r) = +\frac{\omega^4 r^4}{64c^4} E_0 e^{i\omega t}. \quad (23.13)$$

Writing our doubly corrected electric field as $E = E_1 + E_2 + E_3$, we get

$$E = E_0 e^{i\omega t} \left[1 - \frac{1}{2^2} \left(\frac{\omega r}{c} \right)^2 + \frac{1}{2^2 \cdot 4^2} \left(\frac{\omega r}{c} \right)^4 \right]. \quad (23.14)$$

The variation of the electric field with radius is no longer the simple parabola we drew in Fig. 23-5, but at large radii lies slightly above the curve $(E_1 + E_2)$.

We are not quite through yet. The new electric field produces a new correction to the magnetic field, and the newly corrected magnetic field will produce a further correction to the electric field, and on and on. However, we already have all the formulas that we need. For B_3 we can use Eq. (23.10), changing the subscripts of B and E from 2 to 3.

The next correction to the electric field is

$$E_4 = -\frac{1}{2^2 \cdot 4^2 \cdot 6^2} \left(\frac{\omega r}{c} \right)^6 E_0 e^{i\omega t}.$$

So to this order we have that the complete electric field is given by

$$E = E_0 e^{i\omega t} \left[1 - \frac{1}{(1!)^2} \left(\frac{\omega r}{2c} \right)^2 + \frac{1}{(2!)^2} \left(\frac{\omega r}{2c} \right)^4 - \frac{1}{(3!)^2} \left(\frac{\omega r}{2c} \right)^6 + \dots \right], \quad (23.15)$$

where we have written the numerical coefficients in such a way that it is obvious how the series is to be continued.

Our final result is that the electric field between the plates of the capacitor, for any frequency, is given by $E_0 e^{i\omega t}$ times the infinite series which contains only the variable $\omega r/c$. If we wish, we can define a special function, which we will call $J_0(x)$, as the infinite series that appears in the brackets of Eq. (23.15):

$$J_0(x) = 1 - \frac{1}{(1!)^2} \left(\frac{x}{2} \right)^2 + \frac{1}{(2!)^2} \left(\frac{x}{2} \right)^4 - \frac{1}{(3!)^2} \left(\frac{x}{2} \right)^6 + \dots \quad (23.16)$$

Then we can write our solution as $E_0 e^{i\omega t}$ times this function, with $x = \omega r/c$:

$$E = E_0 e^{i\omega t} J_0 \left(\frac{\omega r}{c} \right). \quad (23.17)$$

The reason we have called our special function J_0 is that, naturally, this is not the first time anyone has ever worked out a problem with oscillations in a cylinder. The function has come up before and is usually called J_0 . It always comes up whenever you solve a problem about waves with cylindrical symmetry. The function J_0 is to cylindrical waves what the cosine function is to waves on a straight line. So it is an important function, invented a long time ago. Then a man named Bessel got his name attached to it. The subscript zero means that Bessel invented a whole lot of different functions and this is just the first of them.

The other functions of Bessel— J_1 , J_2 , and so on—have to do with cylindrical waves which have a variation of their strength with the angle around the axis of the cylinder.

The completely corrected electric field between the plates of our circular capacitor, given by Eq. (23.17), is plotted as the solid line in Fig. 23-5. For frequencies that are not too high, our second approximation was already quite good. The third approximation was even better—so good, in fact, that if we had plotted it, you would not have been able to see the difference between it and the solid curve. You will see in the next section, however, that the complete series is needed to get an accurate description for large radii, or for high frequencies.

23-3 A resonant cavity

We want to look now at what our solution gives for the electric field between the plates of the capacitor as we continue to go to higher and higher frequencies. For large ω , the parameter $x = \omega r/c$ also gets large, and the first few terms in the series for J_0 of x will increase rapidly. That means that the parabola we have drawn in Fig. 23-5 curves downward more steeply at higher frequencies. In fact, it looks as though the field would fall all the way to zero at some high frequency, perhaps when c/ω is approximately one-half of a . Let's see whether J_0 does indeed go through zero and become negative. We begin by trying $x = 2$.

$$J_0(2) = 1 - 1 + \frac{1}{4} - \frac{1}{36} = 0.22$$

The function is still not zero, so let's try a higher value of x , say, $x = 2.5$. Putting in numbers, we write

$$J_0(2.5) = 1 - 1.56 + 0.61 - 0.09 = -0.04.$$

The function J_0 has already gone through zero by the time we get to $x = 2.5$. Comparing the results for $x = 2$ and $x = 2.5$, it looks as though J_0 goes through zero at one-fifth of the way from 2.5 to 2. We would guess that the zero occurs for x approximately equal to 2.4. Let's see what that value of x gives:

$$J_0(2.4) = 1 - 1.44 + 0.52 - 0.08 = 0.00$$

We get zero to the accuracy of our two decimal places. If we make the calculation more accurate (or since J_0 is a well-known function, if we look it up in a book), we find that it goes through zero at $x = 2.405$. We have worked it out by hand to show you that you too could have discovered these things rather than having to borrow them from a book.

As long as we are looking up J_0 in a book, it is interesting to notice how it goes for larger values of x , it looks like the graph in Fig. 23-6. As x increases, $J_0(x)$ oscillates between positive and negative values with a decreasing amplitude of oscillation.

We have gotten the following interesting result: If we go high enough in frequency, the electric field at the center of our condenser will be one way and the electric field near the edge will point in the opposite direction. For example,

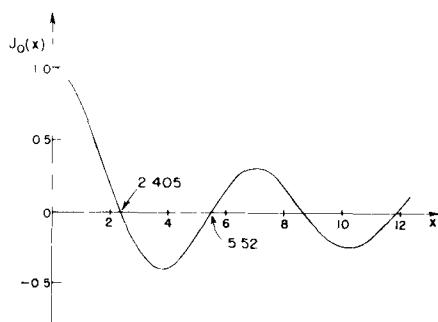


Fig. 23-6. The Bessel function $J_0(x)$.

suppose that we take an ω high enough so that $x = \omega r/c$ at the outer edge of the capacitor is equal to 4; then the edge of the capacitor corresponds to the abscissa $x = 4$ in Fig. 23-6. This means that our capacitor is being operated at the frequency $\omega = 4c/a$. At the edge of the plates, the electric field will have a rather high magnitude opposite the direction we would expect. That is the terrible thing that can happen to a capacitor at high frequencies. If we go to very high frequencies, the direction of the electric field oscillates back and forth many times as we go out from the center of the capacitor. Also there are the magnetic fields associated with these electric fields. It is not surprising that our capacitor doesn't look like the ideal capacitance for high frequencies. We may even start to wonder whether it looks more like a capacitor or an inductance. We should emphasize that there are even more complicated effects that we have neglected which happen at the edges of the capacitor. For instance, there will be a radiation of waves out past the edges, so the fields are even more complicated than the ones we have computed, but we will not worry about those effects now.

We could try to figure out an equivalent circuit for the capacitor, but perhaps it is better if we just admit that the capacitor we have designed for low-frequency fields is just no longer satisfactory when the frequency is too high. If we want to treat the operation of such an object at high frequencies, we should abandon the approximations to Maxwell's equations that we have made for treating circuits and return to the complete set of equations which describe completely the fields in space. Instead of dealing with idealized circuit elements, we have to deal with the real conductors as they are, taking into account all the fields in the spaces in between. For instance, if we want a resonant circuit at high frequencies we will not try to design one using a coil and a parallel-plate capacitor.

We have already mentioned that the parallel-plate capacitor we have been analyzing has some of the aspects of both a capacitor and an inductance. With the electric field there are charges on the surfaces of the plates, and with the magnetic fields there are back emf's. Is it possible that we already have a resonant circuit? We do indeed. Suppose we pick a frequency for which the electric field pattern falls to zero at some radius inside the edge of the disc; that is, we choose $\omega a/c$ greater than 2.405. Everywhere on a circle coaxial with the plates the electric field will be zero. Now suppose we take a thin metal sheet and cut a strip just wide enough to fit between the plates of the capacitor. Then we bend it into a cylinder that will go around at the radius where the electric field is zero. Since there are no electric fields there, when we put this conducting cylinder in place, no currents will flow in it; and there will be no changes in the electric and magnetic fields. We have been able to put a direct short circuit across the capacitor without changing anything. And look what we have; we have a complete cylindrical can with electrical and magnetic fields inside and no connection at all to the outside world. The fields inside won't change even if we throw away the edges of the plates outside our can, and also the capacitor leads. All we have left is a closed can with electric and magnetic fields inside, as shown in Fig. 23-7(a). The electric fields are oscillating back and forth at the frequency ω —which, don't forget, determined the diameter of the can. The amplitude of the oscillating E field varies with the distance from the axis of the can, as shown in the graph of Fig. 23-7(b). This curve is just the first arch of the Bessel function of zero order. There is also a magnetic field which goes in circles around the axis and oscillates in time 90° out of phase with the electric field.

We can also write out a series for the magnetic field and plot it, as shown in the graph of Fig. 23-7(c).

How is it that we can have an electric and magnetic field inside a can with no external connections? It is because the electric and magnetic fields maintain themselves: the changing E makes a B and the changing B makes an E —all according to the equations of Maxwell. The magnetic field has an inductive aspect, and the electric field a capacitive aspect; together they make something like a resonant circuit. Notice that the conditions we have described would only happen if the radius of the can is exactly $2.405 c/\omega$. For a can of a given radius, the oscillating electric and magnetic fields will maintain themselves—in the way we have described

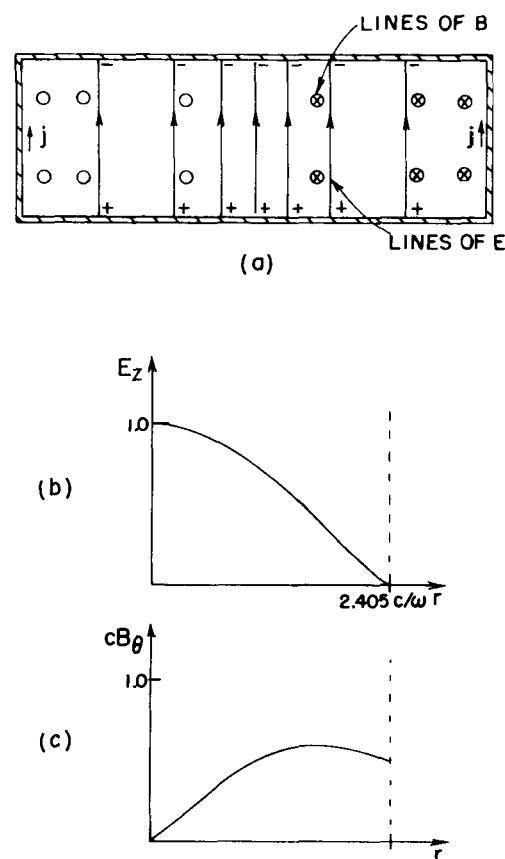


Fig. 23-7. The electric and magnetic fields in an enclosed cylindrical can.

—only at that particular frequency. So a cylindrical can of radius r is *resonant* at the frequency

$$\omega_0 = 2.405 \frac{c}{r}. \quad (23.18)$$

We have said that the fields continue to oscillate in the same way after the can is completely closed. That is not exactly right. It would be possible if the walls of the can were perfect conductors. For a real can, however, the oscillating currents which exist on the inside walls of the can lose energy because of the resistance of the material. The oscillations of the fields will gradually die away. We can see from Fig. 23-7 that there must be strong currents associated with electric and magnetic fields inside the cavity. Because the vertical electrical field stops suddenly at the top and bottom plates of the can, it has a large divergence there; so there must be positive and negative electric charges on the inner surfaces of the can, as shown in Fig. 23-7(a). When the electric field reverses, the charges must reverse also, so there must be an alternating current between the top and bottom plates of the can. These charges will flow in the sides of the can, as shown in the figure. We can also see that there must be currents in the sides of the can by considering what happens to the magnetic field. The graph of Fig. 23-7(c) tells us that the magnetic field suddenly drops to zero at the edge of the can. Such a sudden change in the magnetic field can happen only if there is a current in the wall. This current is what gives the alternating electric charges on the top and bottom plates of the can.

You may be wondering about our discovery of currents in the vertical sides of the can. What about our earlier statement that nothing would be changed when we introduced these vertical sides in a region where the electric field was zero? Remember, however, that when we first put in the sides of the can, the top and bottom plates extended out beyond them, so that there were also magnetic fields on the outside of our can. It was only when we threw away the parts of the capacitor plates beyond the edges of the can that net currents had to appear on the insides of the vertical walls.

Although the electric and magnetic fields in the completely enclosed can will gradually die away because of the energy losses, we can stop this from happening if we make a little hole in the can and put in a little bit of electrical energy to make up the losses. We take a small wire, poke it through the hole in the side of the can, and fasten it to the inside wall so that it makes a small loop, as shown in Fig. 23-8. If we now connect this wire to a source of high-frequency alternating current, this current will couple energy into the electric and magnetic fields of the cavity and keep the oscillations going. This will happen, of course, only if the frequency of the driving source is at the resonant frequency of the can. If the source is at the wrong frequency, the electric and magnetic fields will not resonate, and the fields in the can will be very weak.

The resonant behavior can easily be seen by making another small hole in the can and hooking in another coupling loop, as we have also drawn in Fig. 23-8. The changing magnetic field through this loop will generate an induced electromotive force in the loop. If this loop is now connected to some external measuring circuit, the currents will be proportional to the strength of the fields in the cavity. Suppose we now connect the input loop of our cavity to an RF signal generator, as shown in Fig. 23-9. The signal generator contains a source of alternating current whose frequency can be varied by varying the knob on the front of the generator. Then we connect the output loop of the cavity to a “detector,” which is an instrument that measures the current from the output loop. It gives a meter reading proportional to this current. If we now measure the output current as a function of the frequency of the signal generator, we find a curve like that shown in Fig. 23-10. The output current is small for all frequencies except those very near the frequency ω_0 , which is the resonant frequency of the cavity. The resonance curve is very much like those we described in Chapter 23 of Vol. I. The width of the resonance is, however, much narrower than we usually find for resonant circuits made of inductances and capacitors; that is, the Q of the cavity is very high. It is not unusual to find Q 's as high as 100,000 or more if the inside walls of the cavity are made of some material with a very good conductivity, such as silver.

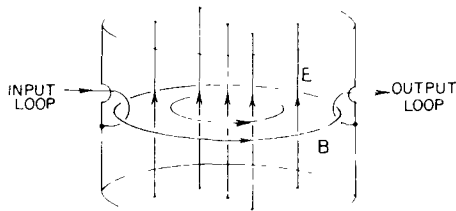


Fig. 23-8. Coupling into and out of a resonant cavity.

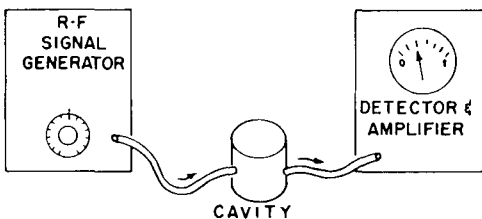


Fig. 23-9. A setup for observing the cavity resonance.

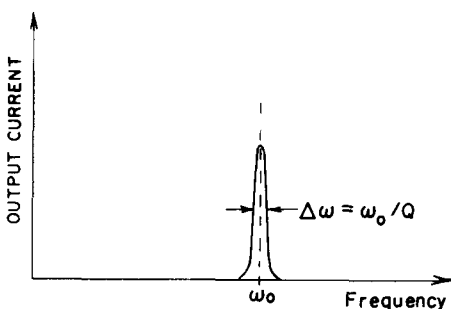


Fig. 23-10. The frequency response curve of a resonant cavity.

23-4 Cavity modes

Suppose we now try to check our theory by making measurements with an actual can. We take a can which is a cylinder with a diameter of 3.0 inches and a height of about 2.5 inches. The can is fitted with an input and output loop, as shown in Fig. 23-8. If we calculate the resonant frequency expected for this can according to Eq (23.18), we get that $f_0 = \omega_0/2\pi = 3010$ megacycles. When we set the frequency of our signal generator near 3000 megacycles and vary it slightly until we find the resonance, we observe that the maximum output current occurs for a frequency of 3050 megacycles, which is quite close to the predicted resonant frequency, but not exactly the same. There are several possible reasons for the discrepancy. Perhaps the resonant frequency is changed a little bit because of the holes we have cut to put in the coupling loops. A little thought, however, shows that the holes should lower the resonant frequency a little bit, so that cannot be the reason. Perhaps there is some slight error in the frequency calibration of the signal generator, or perhaps our measurement of the diameter of the cavity is not accurate enough. Anyway, the agreement is fairly close.

Much more important is something that happens if we vary the frequency of our signal generator somewhat further from 3000 megacycles. When we do that we get the results shown in Fig. 23-11. We find that, in addition to the resonance we expected near 3000 megacycles, there is also a resonance near 3300 megacycles and one near 3820 megacycles. What do these extra resonances mean? We might get a clue from Fig. 23-6. Although we have been assuming that the first zero of the Bessel function occurs at the edge of the can, it could also be that the second zero of the Bessel function corresponds to the edge of the can, so that there is one complete oscillation of the electric field as we move from the center of the can out to the edge, as shown in Fig. 23-12. This is another possible mode for the oscillating fields. We should certainly expect the can to resonate in such a mode. But notice, the second zero of the Bessel function occurs at $x = 5.52$, which is over twice as large as the value at the first zero. The resonant frequency of this mode should therefore be higher than 6000 megacycles. We would, no doubt, find it there, but it doesn't explain the resonance we observe at 3300.

The trouble is that in our analysis of the behavior of a resonant cavity we have considered only one possible geometric arrangement of the electric and magnetic fields. We have assumed that the electric fields are vertical and that the magnetic fields lie in horizontal circles. But other fields are possible. The only requirements are that the fields should satisfy Maxwell's equations inside the can and that the electric field should meet the wall at right angles. We have considered the case in which the top and the bottom of the can are flat, but things would not be completely different if the top and bottom were curved. In fact, how is the can supposed to know which is its top and bottom, and which are its sides? It is, in fact, possible to show that there is a mode of oscillation of the fields inside the can in which the electric fields go more or less across the diameter of the can, as shown in Fig. 23-13.

It is not too hard to understand why the natural frequency of this mode should be not very different from the natural frequency of the first mode we have considered. Suppose that instead of our cylindrical cavity we had taken a cavity which was a cube 3 inches on a side. It is clear that this cavity would have three different modes, but all with the same frequency. A mode with the electric field going more or less up and down would certainly have the same frequency as the mode in which the electric field was directed right and left. If we now distort the cube into a cylinder, we will change these frequencies somewhat. We would still expect them not to be changed too much, provided we keep the dimensions of the cavity more or less the same. So the frequency of the mode of Fig. 23-13 should not be too different from the mode of Fig. 23-8. We could make a detailed calculation of the natural frequency of the mode shown in Fig. 23-13, but we will not do that now. When the calculations are carried through, it is found that, for the dimensions we have assumed, the resonant frequency comes out very close to the observed resonance at 3300 megacycles.

By similar calculations it is possible to show that there should be still another mode at the other resonant frequency we found near 3800 megacycles. For this

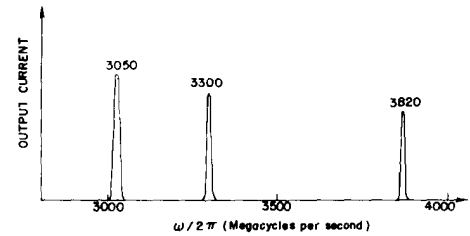
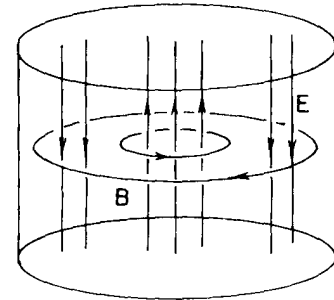
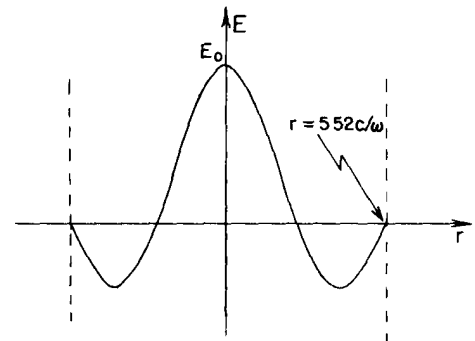


Fig. 23-11. Observed resonant frequencies of a cylindrical cavity.



(a)



(b)

Fig. 23-12 A higher-frequency mode.

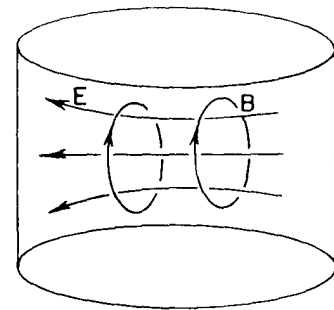


Fig. 23-13. A transverse mode of the cylindrical cavity.

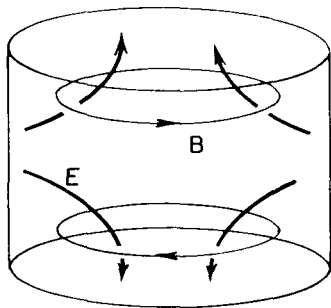


Fig. 23-14. Another mode of a cylindrical cavity.

mode, the electric and magnetic fields are as shown in Fig. 23-14. The electric field does not bother to go all the way across the cavity. It goes from the sides to the ends, as shown.

As you will probably now believe, if we go higher and higher in frequency we should expect to find more and more resonances. There are many different modes, each of which will have a different resonant frequency corresponding to some particular complicated arrangement of the electric and magnetic fields. Each of these field arrangements is called a resonant *mode*. The resonance frequency of each mode can be calculated by solving Maxwell's equations for the electric and magnetic fields in the cavity.

When we have a resonance at some particular frequency, how can we know which mode is being excited? One way is to poke a little wire into the cavity through a small hole. If the electric field is along the wire, as in Fig. 23-15(a), there will be relatively large currents in the wire, sapping energy from the fields, and the resonance will be suppressed. If the electric field is as shown in Fig. 23-15(b), the wire will have a much smaller effect. We could find which way the field points in this mode by bending the end of the wire, as shown in Fig. 23-15(c). Then, as we rotate the wire, there will be a big effect when the end of the wire is parallel to E and a small effect when it is rotated so as to be at 90° to E .

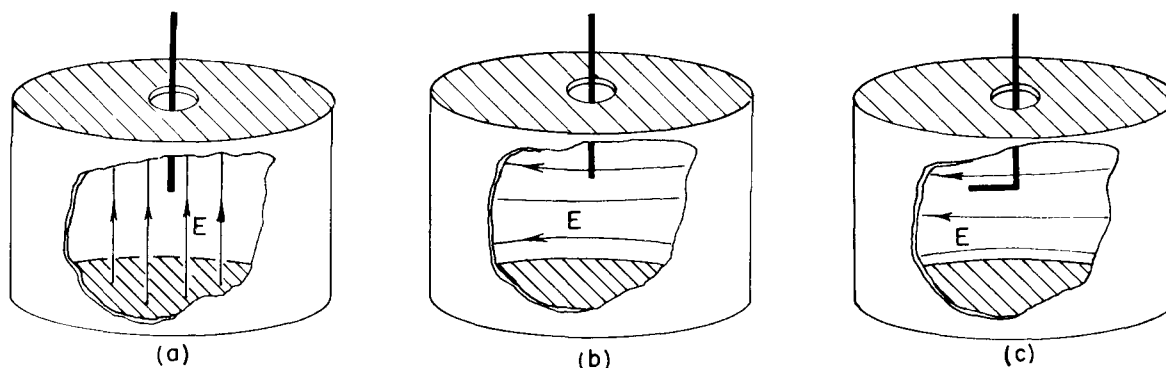


Fig. 23-15. A short metal wire inserted into a cavity will disturb the resonance much more when it is parallel to E than when it is at right angles.

23-5 Cavities and resonant circuits

Although the resonant cavity we have been describing seems to be quite different from the ordinary resonant circuit consisting of an inductance and a capacitor, the two resonant systems are, of course, closely related. They are both members of the same family; they are just two extreme cases of electromagnetic resonators—and there are many intermediate cases between these two extremes. Suppose we start by considering the resonant circuit of a capacitor in parallel with an inductance, as shown in Fig. 23-16(a). This circuit will resonate at the frequency $\omega_0 = 1/\sqrt{LC}$. If we want to raise the resonant frequency of this circuit, we can do so by lowering the inductance L . One way is to decrease the number of turns in the coil. We can, however, go only so far in this direction. Eventually we will get down to the last turn, and we will have just a piece of wire joining the top and bottom plates of the condenser. We could raise the resonant frequency still further by making the capacitance smaller; however, we can also continue to decrease the inductance by putting several inductances in parallel. Two one-turn inductances in parallel will have only half the inductance of each turn. So when our inductance has been reduced to a single turn, we can continue to raise the resonant frequency by adding other single loops from the top plate to the bottom plate of the condenser. For instance, Fig. 23-16(b) shows the condenser plates connected by six such "single-turn inductances." If we continue to add many such pieces of wire, we can make the transition to the completely enclosed resonant system shown in part (c) of the figure, which is a drawing of the cross section of a cylindrically symmetrical

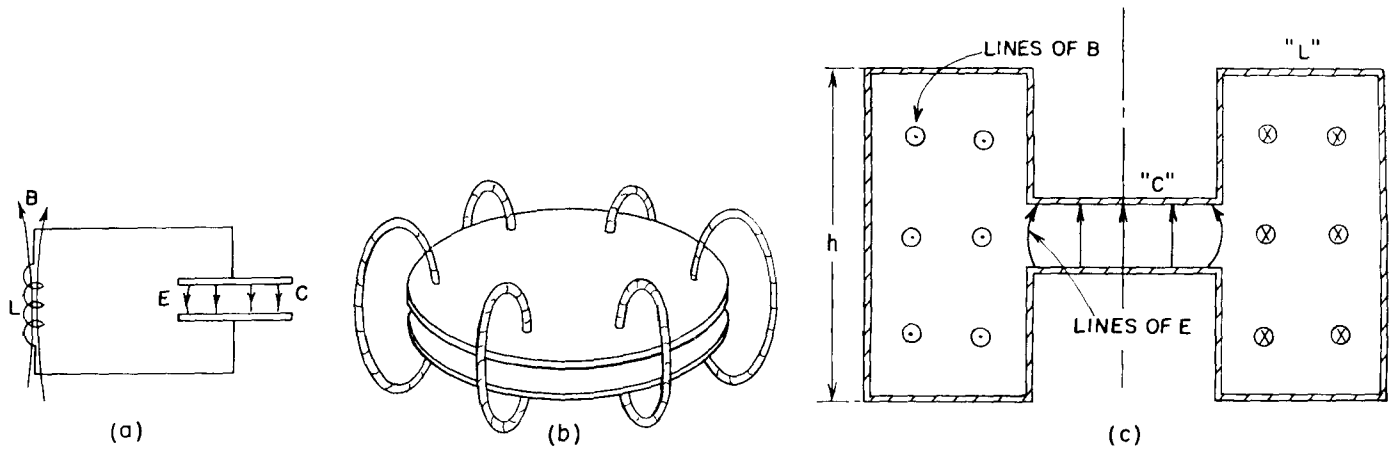


Fig. 23-16. Resonators of progressively higher resonant frequencies.

object. Our inductance is now a cylindrical hollow can attached to the edges of the condenser plates. The electric and magnetic fields will be as shown in the figure. Such an object is, of course, a resonant cavity. It is called a “loaded” cavity. But we can still think of it as an L - C circuit in which the capacity section is the region where we find most of the electric field and the inductance section is that region where we find most of the magnetic field.

If we want to make the frequency of the resonator in Fig. 23-16(c) still higher, we can do so by continuing to decrease the inductance L . To do that, we must decrease the geometric dimensions of the inductance section, for example by decreasing the dimension h in the drawing. As h is decreased, the resonant frequency will be increased. Eventually, of course, we will get to the situation in which the height h is just equal to the separation between the condenser plates. We then have just a cylindrical can, our resonant circuit has become the cavity resonator of Fig. 23-7.

You will notice that in the original L - C resonant circuit of Fig. 23-16 the electric and magnetic fields are quite separate. As we have gradually modified the resonant system to make higher and higher frequencies, the magnetic field has been brought closer and closer to the electric field until in the cavity resonator the two are quite intermixed.

Although the cavity resonators we have talked about in this chapter have been cylindrical cans, there is nothing magic about the cylindrical shape. A can of any shape will have resonant frequencies corresponding to various possible modes of oscillations of the electric and magnetic fields. For example, the “cavity” shown in Fig. 23-17 will have its own particular set of resonant frequencies—although they would be rather difficult to calculate.

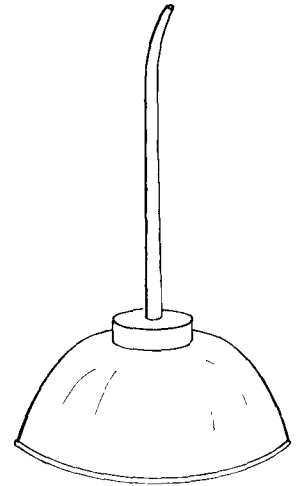


Fig. 23-17. Another resonant cavity.

Waveguides

24-1 The transmission line

In the last chapter we studied what happened to the lumped elements of circuits when they were operated at very high frequencies, and we were led to see that a resonant circuit could be replaced by a cavity with the fields resonating inside. Another interesting technical problem is the connection of one object to another, so that electromagnetic energy can be transmitted between them. In low-frequency circuits the connection is made with wires, but this method doesn't work very well at high frequencies because the circuits would radiate energy into all the space around them, and it is hard to control where the energy will go. The fields spread out around the wires; the currents and voltages are not "guided" very well by the wires. In this chapter we want to look into the ways that objects can be interconnected at high frequencies. At least, that's one way of presenting our subject.

Another way is to say that we have been discussing the behavior of waves in free space. Now it is time to see what happens when oscillating fields are confined in one or more dimensions. We will discover the interesting new phenomenon when the fields are confined in only two dimensions and allowed to go free in the third dimension, they propagate in waves. These are "guided waves"—the subject of this chapter.

We begin by working out the general theory of the *transmission line*. The ordinary power transmission line that runs from tower to tower over the countryside radiates away some of its power, but the power frequencies (50–60 cycles/sec) are so low that this loss is not serious. The radiation could be stopped by surrounding the line with a metal pipe, but this method would not be practical for power lines because the voltages and currents used would require a very large, expensive, and heavy pipe. So simple "open lines" are used.

For somewhat higher frequencies—say a few kilocycles—radiation can already be serious. However, it can be reduced by using "twisted-pair" transmission lines, as is done for short-run telephone connections. At higher frequencies, however, the radiation soon becomes intolerable, either because of power losses or because the energy appears in other circuits where it isn't wanted. For frequencies from a few kilocycles to some hundreds of megacycles, electromagnetic signals and power are usually transmitted via coaxial lines consisting of a wire inside a cylindrical "outer conductor" or "shield." Although the following treatment will apply to a transmission line of two parallel conductors of any shape, we will carry it out referring to a coaxial line.

We take the simplest coaxial line that has a central conductor, which we suppose is a thin hollow cylinder, and an outer conductor which is another thin cylinder on the same axis as the inner conductor, as in Fig. 24-1. We begin by figuring out approximately how the line behaves at relatively low frequencies. We have already described some of the low-frequency behavior when we said earlier that two such conductors had a certain amount of inductance per unit length or a certain capacity per unit length. We can, in fact, describe the low-frequency behavior of any transmission line by giving its inductance per unit length, L_0 and its capacity per unit length, C_0 . Then we can analyze the line as the limiting case of the L - C filter as discussed in Section 22-6. We can make a filter which imitates the line by taking small series elements $L_0 \Delta x$ and small shunt capacities $C_0 \Delta x$, where Δx is an element of length of the line. Using our results for the infinite filter, we see that there would be a propagation of electric

24-1 The transmission line

24-2 The rectangular waveguide

24-3 The cutoff frequency

24-4 The speed of the guided waves

24-5 Observing guided waves

24-6 Waveguide plumbing

24-7 Waveguide modes

24-8 Another way of looking at the guided waves

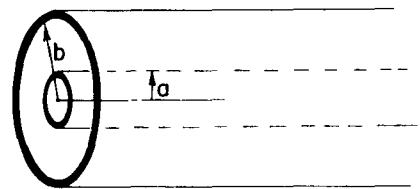


Fig. 24-1. A coaxial transmission line.

signals along the line. Rather than following that approach, however, we would now rather look at the line from the point of view of a differential equation.

Suppose that we see what happens at two neighboring points along the transmission line, say at the distances x and $x + \Delta x$ from the beginning of the line. Let's call the voltage difference between the two conductors $V(x)$, and the current along the "hot" conductor $I(x)$ (see Fig. 24-2). If the current in the line is varying, the inductance will give us a voltage drop across the small section of line from x to $x + \Delta x$ in the amount

$$\Delta V = V(x + \Delta x) - V(x) = -L_0 \Delta x \frac{dI}{dt}.$$

Or, taking the limit as $\Delta x \rightarrow 0$, we get

$$\frac{\partial V}{\partial x} = -L_0 \frac{\partial I}{\partial t}. \quad (24.1)$$

The changing current gives a gradient of the voltage.

Referring again to the figure, if the voltage at x is changing, there must be some charge supplied to the capacity in that region. If we take the small piece of line between x and $x + \Delta x$, the charge on it is $q = C_0 \Delta x V$. The time rate-of-change of this charge is $C_0 \Delta x dV/dt$, but the charge changes only if the current $I(x)$ into the element is different from the current $I(x + \Delta x)$ out. Calling the difference ΔI , we have

$$\Delta I = -C_0 \Delta x \frac{dV}{dt}.$$

Taking the limit as $\Delta x \rightarrow 0$, we get

$$\frac{\partial I}{\partial x} = -C_0 \frac{\partial V}{\partial t}. \quad (24.2)$$

So the conservation of charge implies that the gradient of the current is proportional to the time rate-of-change of the voltage.

Equations (24.1) and (24.2) are then the basic equations of a transmission line. If we wish, we could modify them to include the effects of resistance in the conductors or of leakage of charge through the insulation between the conductors, but for our present discussion we will just stay with the simple example.

The two transmission line equations can be combined by differentiating one with respect to t and the other with respect to x and eliminating either V or I . Then we have either

$$\frac{\partial^2 V}{\partial x^2} = C_0 L_0 \frac{\partial^2 V}{\partial t^2} \quad (24.3)$$

or

$$\frac{\partial^2 I}{\partial x^2} = C_0 L_0 \frac{\partial^2 I}{\partial t^2}. \quad (24.4)$$

Once more we recognize the wave equation in x . For a uniform transmission line, the voltage (and current) propagates along the line as a wave. The voltage along the line must be of the form $V(x, t) = f(x - vt)$ or $V(x, t) = g(x + vt)$, or a sum of both. Now what is the velocity v ? We know that the coefficient of the $\partial^2/\partial t^2$ term is just $1/v^2$, so

$$v = \frac{1}{\sqrt{L_0 C_0}}. \quad (24.5)$$

We will leave it for you to show that the voltage for *each wave* in a line is proportional to the current of that wave and that the constant of proportionality is just the characteristic impedance z_0 . Calling V_+ and I_+ the voltage and current for a wave going in the plus x -direction, you should get

$$V_+ = z_0 I_+. \quad (24.6)$$

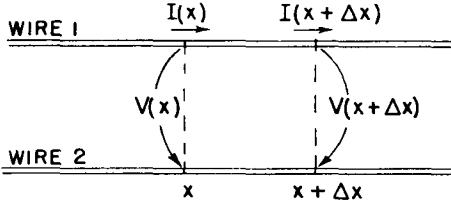


Fig. 24-2. The currents and voltages of a transmission line.

Similary, for the wave going toward minus x the relation is

$$V_- = -z_0 I_-.$$

The characteristic impedance—as we found out from our filter equations—is given by

$$z_0 = \sqrt{\frac{L_0}{C_0}}, \quad (24.7)$$

and is, therefore, a pure resistance.

To find the propagation speed v and the characteristic impedance z_0 of a transmission line, we have to know the inductance and capacity per unit length. We can calculate them easily for a coaxial cable, so we will see how that goes. For the inductance we follow the ideas of Section 17-8, and set $\frac{1}{2}LI^2$ equal to the magnetic energy which we get by integrating $\epsilon_0 c^2 B^2/2$ over the volume. Suppose that the central conductor carries the current I ; then we know that $B = I/2\pi\epsilon_0 c^2 r$, where r is the distance from the axis. Taking as a volume element a cylindrical shell of thickness dr and of length l , we have for the magnetic energy

$$U = \frac{\epsilon_0 c^2}{2} \int_a^b \left(\frac{I}{2\pi\epsilon_0 c^2 r} \right)^2 l 2\pi r dr,$$

where a and b are the radii of the inner and outer conductors, respectively. Carrying out the integral, we get

$$U = \frac{I^2 l}{4\pi\epsilon_0 c^2} \ln \frac{b}{a}. \quad (24.8)$$

Setting the energy equal to $\frac{1}{2}LI^2$, we find

$$L = \frac{l}{2\pi\epsilon_0 c^2} \ln \frac{b}{a}. \quad (24.9)$$

It is, as it should be, proportional to the length l of the line, so the inductance per unit length L_0 is

$$L_0 = \frac{\ln(b/a)}{2\pi\epsilon_0 c^2}. \quad (24.10)$$

We have worked out the charge on a cylindrical condenser (see Section 12-2). Now, dividing the charge by the potential difference, we get

$$C = \frac{2\pi\epsilon_0 l}{\ln(b/a)}.$$

The capacity per unit length C_0 is C/l . Combining this result with Eq. (24.10), we see that the product $L_0 C_0$ is just equal to $1/c^2$, so $v = 1/\sqrt{L_0 C_0}$ is equal to c . The wave travels down the line with the speed of light. We point out that this result depends on our assumptions: (a) that there are no dielectrics or magnetic materials in the space between the conductors, and (b) that the currents are all on the surfaces of the conductors (as they would be for perfect conductors). We will see later that for good conductors at high frequencies, all currents distribute themselves on the surfaces as they would for a perfect conductor, so this assumption is then valid.

Now it is interesting that so long as assumptions (a) and (b) are correct, the product $L_0 C_0$ is equal to $1/c^2$ for *any* parallel pair of conductors—even, say, for a hexagonal inner conductor anywhere inside an elliptical outer conductor. So long as the cross section is constant and the space between has no material, waves are propagated at the velocity of light.

No such general statement can be made about the characteristic impedance. For the coaxial line, it is

$$z_0 = \frac{\ln(b/a)}{2\pi\epsilon_0 c}. \quad (24.11)$$

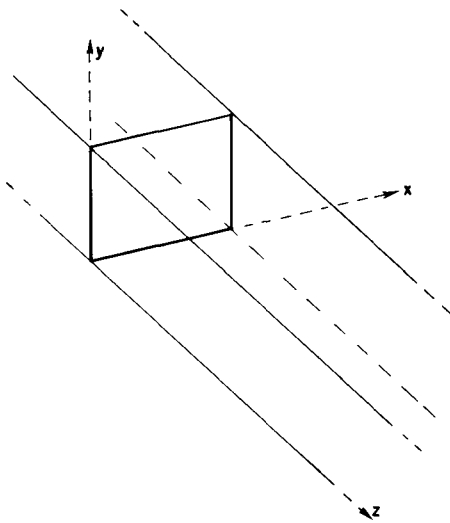


Fig. 24-3. Coordinates chosen for the rectangular waveguide.

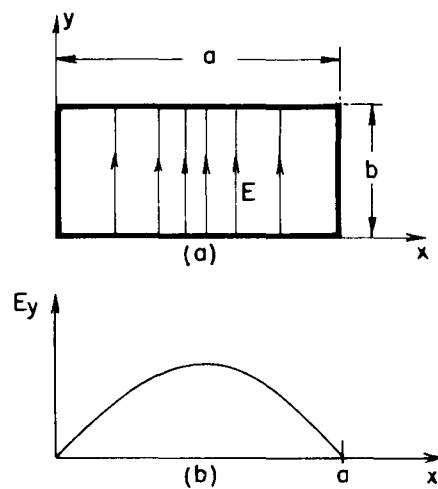


Fig. 24-4. The electric field in the waveguide at some value of z .

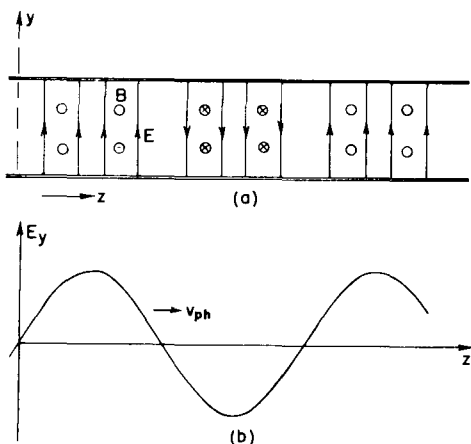


Fig. 24-5. The z -dependence of the field in the waveguide.

The factor $1/\epsilon_0 c$ has the dimensions of a resistance and is equal to 120π ohms. The geometric factor $\ln(b/a)$ depends only logarithmically on the dimensions, so for the coaxial line—and most lines—the characteristic impedance has typical values of from 50 ohms or so to a few hundred ohms.

24-2 The rectangular waveguide

The next thing we want to talk about seems, at first sight, to be a striking phenomenon: if the central conductor is removed from the coaxial line, it can still carry electromagnetic power. In other words, at high enough frequencies a hollow tube will work just as well as one with wires. It is related to the mysterious way in which a resonant circuit of a condenser and inductance gets replaced by nothing but a can at high frequencies.

Although it may seem to be a remarkable thing when one has been thinking in terms of a transmission line as a distributed inductance and capacity, we all know that electromagnetic waves can travel along inside a hollow metal pipe. If the pipe is straight, we can see through it! So certainly electromagnetic waves go through a pipe. But we also know that it is not possible to transmit low-frequency waves (power or telephone) through the inside of a single metal pipe. So it must be that electromagnetic waves will go through if their wavelength is short enough. Therefore we want to discuss the limiting case of the longest wavelength (or the lowest frequency) that can get through a pipe of a given size. Since the pipe is then being used to carry waves, it is called a *waveguide*.

We will begin with a rectangular pipe, because it is the simplest case to analyze. We will first give a mathematical treatment and come back later to look at the problem in a much more elementary way. The more elementary approach, however, can be applied easily only to a rectangular guide. The basic phenomena are the same for a general guide of arbitrary shape, so the mathematical argument is fundamentally more sound.

Our problem, then, is to find what kind of waves can exist inside a rectangular pipe. Let's first choose some convenient coordinates; we take the z -axis along the length of the pipe, and the x - and y -axes parallel to the two sides, as shown in Fig. 24-3.

We know that when light waves go down the pipe, they have a transverse electric field; so suppose we look first for solutions in which \mathbf{E} is perpendicular to z , say with only a y -component, E_y . This electric field will have some variation across the guide; in fact, it must go to zero at the sides parallel to the y -axis, because the currents and charges in a conductor always adjust themselves so that there is no tangential component of the electric field at the surface of a conductor. So E_y will vary with x in some arch, as shown in Fig. 24-4. Perhaps it is the Bessel function we found for a cavity? No, because the Bessel function has to do with cylindrical geometries. For a rectangular geometry, waves are usually simple harmonic functions, so we should try something like $\sin k_x x$.

Since we want waves that propagate down the guide, we expect the field to alternate between positive and negative values as we go along in z , as in Fig. 24-5, and these oscillations will travel along the guide with some velocity v . If we have oscillations at some definite frequency ω , we would guess that the wave might vary with z like $\cos(\omega t - k_z z)$, or to use the more convenient mathematical form, like $e^{i(\omega t - k_z z)}$. This z -dependence represents a wave travelling with the speed $v = \omega/k_z$ (see Chapter 29, Vol. I).

So we might guess that the wave in the guide would have the following mathematical form:

$$E_y = E_0 \sin k_x x e^{i(\omega t - k_z z)}. \quad (24.12)$$

Let's see whether this guess satisfies the correct field equations. First, the electric field should have no tangential components at the conductors. Our field satisfies this requirement; it is perpendicular to the top and bottom faces and is zero at the two side faces. Well, it is if we choose k_x so that one-half a cycle of

$\sin k_x x$ just fits in the width of the guide—that is, if

$$k_x a = \pi. \quad (24.13)$$

There are other possibilities, like $k_x a = 2\pi, 3\pi, \dots$, or, in general,

$$k_x a = n\pi, \quad (24.14)$$

where n is any integer. These represent various complicated arrangements of the field, but for now let's take only the simplest one, where $k_x = \pi/a$, where a is the width of the inside of the guide.

Next, the divergence of \mathbf{E} must be zero in the free space inside the guide, since there are no charges there. Our \mathbf{E} has only a y -component, and it doesn't change with y , so we do have that $\nabla \cdot \mathbf{E} = 0$.

Finally, our electric field must agree with the rest of Maxwell's equations in the free space inside the guide. That is the same thing as saying that it must satisfy the wave equation

$$\frac{\partial^2 E_y}{\partial x^2} + \frac{\partial^2 E_y}{\partial y^2} + \frac{\partial^2 E_y}{\partial z^2} - \frac{1}{c^2} \frac{\partial^2 E_y}{\partial t^2} = 0. \quad (24.15)$$

We have to see whether our guess, Eq. (24.12), will work. The second derivative of E_y with respect to x is just $-k_x^2 E_y$. The second derivative with respect to y is zero, since nothing depends on y . The second derivative with respect to z is $-k_z^2 E_y$, and the second derivative with respect to t is $-\omega^2 E_y$. Equation (24.15) then says that

$$k_x^2 E_y + k_z^2 E_y - \frac{\omega^2}{c^2} E_y = 0.$$

Unless E_y is zero everywhere (which is not very interesting), this equation is correct if

$$k_x^2 + k_z^2 - \frac{\omega^2}{c^2} = 0. \quad (24.16)$$

We have already fixed k_x , so this equation tells us that there can be waves of the type we have assumed if k_z is related to the frequency ω so that Eq. (24.16) is satisfied—in other words, if

$$k_z = \sqrt{(\omega^2/c^2) - (\pi^2/a^2)}. \quad (24.17)$$

The waves we have described are propagated in the z -direction with this value of k_z .

The wave number k_z we get from Eq. (24.17) tells us, for a given frequency ω , the speed with which the nodes of the wave propagate down the guide. The phase velocity is

$$v = \frac{\omega}{k_z}. \quad (24.18)$$

You will remember that the wavelength λ of a travelling wave is given by $\lambda = 2\pi v/\omega$, so k_z is also equal to $2\pi/\lambda_g$, where λ_g is the wavelength of the oscillations along the z -direction—the “guide wavelength.” The wavelength in the guide is different, of course, from the free-space wavelength of electromagnetic waves of the same frequency. If we call the free-space wavelength λ_0 , which is equal to $2\pi c/\omega$, we can write Eq. (24.17) as

$$\lambda_g = \frac{\lambda_0}{\sqrt{1 - (\lambda_0/2a)^2}}. \quad (24.19)$$

Besides the electric fields there are magnetic fields that will travel with the wave, but we will not bother to work out an expression for them right now. Since $c^2 \nabla \times \mathbf{B} = \partial \mathbf{E} / \partial t$, the lines of \mathbf{B} will circulate around the regions in which $\partial \mathbf{E} / \partial t$ is largest, that is, halfway between the maximum and minimum of \mathbf{E} . The loops of \mathbf{B} will lie parallel to the xz -plane and between the crests and troughs of \mathbf{E} , as shown in Fig. 24-6.

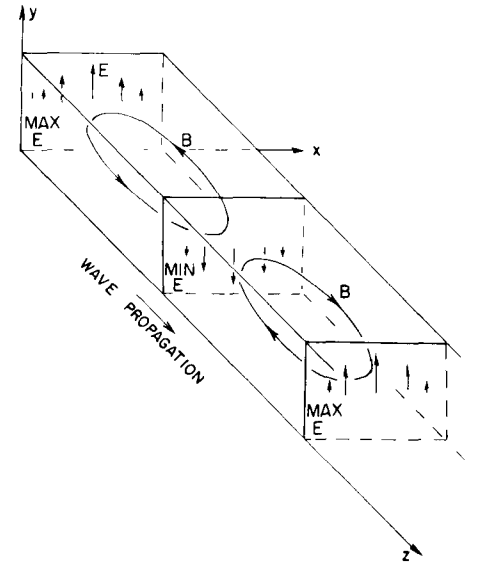


Fig. 24-6. The magnetic field in the waveguide.

24-3 The cutoff frequency

In solving Eq. (24.16) for k_z , there should really be two roots—one plus and one minus. We should write

$$k_z = \pm \sqrt{(\omega^2/c^2) - (\pi^2/a^2)}. \quad (24.20)$$

The two signs simply mean that there can be waves which propagate with a negative phase velocity (toward $-z$), as well as waves which propagate in the positive direction in the guide. Naturally, it should be possible for waves to go in either direction. Since both types of waves can be present at the same time, there will be the possibility of standing-wave solutions.

Our equation for k_z also tells us that higher frequencies give larger values of k_z , and therefore smaller wavelengths, until in the limit of large ω , k becomes equal to ω/c , which is the value we would expect for waves in free space. The light we “see” through a pipe still travels at the speed c . But now notice that if we go toward low frequencies, something strange happens. At first the wavelength gets longer and longer, but if ω gets too small the quantity inside the square root of Eq. (24.20) suddenly becomes negative. This will happen as soon as ω gets to be less than $\pi c/a$ —or when λ_0 becomes greater than $2a$. In other words, when the frequency gets smaller than a certain critical frequency $\omega_c = \pi c/a$, the wave number k_z (and also λ_0) becomes imaginary and we haven’t got a solution any more. Or do we? Who said that k_z has to be real? What if it does come out imaginary? Our field equations are still satisfied. Perhaps an imaginary k_z also represents a wave.

Suppose ω is less than ω_c ; then we can write

$$k_z = \pm ik', \quad (24.21)$$

where k' is a positive real number:

$$k' = \sqrt{(\pi^2/a^2) - (\omega^2/c^2)}. \quad (24.22)$$

If we now go back to our expression, Eq. (24.12), for E_y , we have

$$E_y = E_0 \sin k_x x e^{i(\omega t \mp k' z)}, \quad (24.23)$$

which we can write as

$$E_y = E_0 \sin k_x x e^{\pm k' z} e^{i\omega t}. \quad (24.24)$$

This expression gives an E -field that oscillates with time as $e^{i\omega t}$ but which varies with z as $e^{\pm k' z}$. It decreases or increases with z smoothly as a real exponential. In our derivation we didn’t worry about the sources that started the waves, but there must, of course, be a source someplace in the guide. The sign that goes with k' must be the one that makes the field decrease with increasing distance from the source of the waves.

So for frequencies below $\omega_c = \pi c/a$, waves do *not* propagate down the guide; the oscillating fields penetrate into the guide only a distance of the order of $1/k'$. For this reason, the frequency ω_c is called the “cutoff frequency” of the guide. Looking at Eq. (24.22), we see that for frequencies just a little below ω_c , the number k' is small and the fields can penetrate a long distance into the guide. But if ω is much less than ω_c , the exponential coefficient k' is equal to π/a and the field dies off extremely rapidly, as shown in Fig. 24-7. The field decreases by $1/e$ in the distance a/π , or in only about one-third of the guide width. The fields penetrate very little distance from the source.

We want to emphasize an interesting feature of our analysis of the guided waves—the appearance of the imaginary wave number k_z . Normally, if we solve an equation in physics and get an imaginary number, it doesn’t mean anything physical. For *waves*, however, an imaginary wave number *does* mean something. The wave equation is still satisfied; it only means that the solution gives exponentially decreasing fields instead of propagating waves. So in any wave problem where k becomes imaginary for some frequency, it means that the form of the wave changes—the sine wave changes into an exponential.

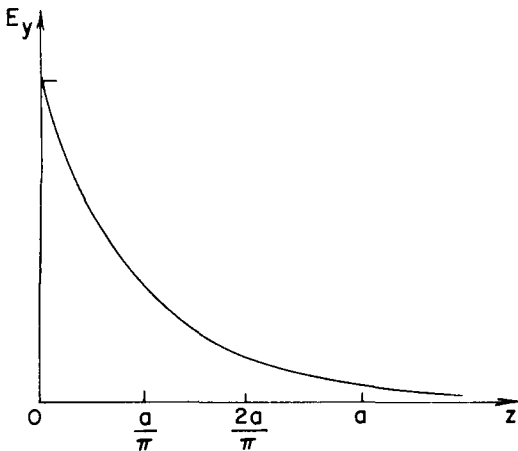


Fig. 24-7. The variation of E_y with z for $\omega \ll \omega_c$.

24-4 The speed of the guided waves

The wave velocity we have used above is the phase velocity, which is the speed of a node of the wave; it is a function of frequency. If we combine Eqs. (24.17) and (24.18), we can write

$$v_{\text{phase}} = \frac{c}{\sqrt{1 - (\omega_c/\omega)^2}}. \quad (24.25)$$

For frequencies above cutoff—where travelling waves exist— ω_c/ω is less than one, and v_{phase} is real and *greater than* the speed of light. We have already seen in Chapter 48 of Vol. I that *phase* velocities greater than light are possible, because it is just the nodes of the wave which are moving and not energy or information. In order to know how fast *signals* will travel, we have to calculate the speed of pulses or modulations made by the interference of a wave of one frequency with one or more waves of slightly different frequencies (see Chapter 48, Vol. I). We have called the speed of the envelope of such a group of waves the group velocity; it is not ω/k but $d\omega/dk$:

$$v_{\text{group}} = \frac{d\omega}{dk}. \quad (24.26)$$

Taking the derivative of Eq. (24.17) with respect to ω and inverting to get $d\omega/dk$, we find that

$$v_{\text{group}} = c \sqrt{1 - (\omega_c/\omega)^2}, \quad (24.27)$$

which is less than the speed of light.

The geometric mean of v_{phase} and v_{group} is just c , the speed of light:

$$v_{\text{phase}} v_{\text{group}} = c^2. \quad (24.28)$$

This is curious, because we have seen a similar relation in quantum mechanics. For a particle with any velocity—even relativistic—the momentum p and energy U are related by

$$U^2 = p^2 c^2 + m^2 c^4. \quad (24.29)$$

But in quantum mechanics the energy is $\hbar\omega$, and the momentum is \hbar/λ , which is equal to $\hbar k$; so Eq. (24.29) can be written

$$\frac{\omega^2}{c^2} = k^2 + \frac{m^2 c^2}{\hbar^2}, \quad (24.30)$$

or

$$k = \sqrt{(\omega^2/c^2) - (m^2 c^2/\hbar^2)}, \quad (24.31)$$

which looks very much like Eq. (24.17) . . . Interesting!

The group velocity of the waves is also the speed at which energy is transported along the guide. If we want to find the energy flow down the guide, we can get it from the energy density times the group velocity. If the root mean square electric field is E_0 , then the average density of electric energy is $\epsilon_0 E_0^2/2$. There is also some energy associated with the magnetic field. We will not prove it here, but in any cavity or guide the magnetic and electric energies are equal, so the total electromagnetic energy density is $\epsilon_0 E_0^2$. The power dU/dt transmitted by the guide is then

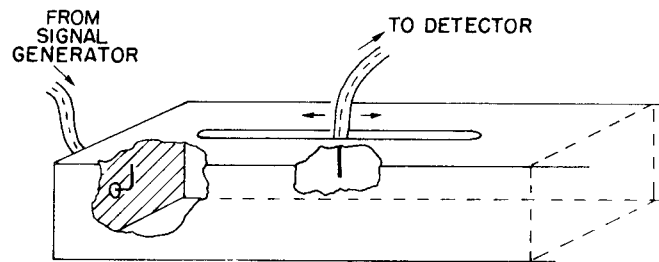
$$\frac{dU}{dt} = \epsilon_0 E_0^2 a b v_{\text{group}}. \quad (24.32)$$

(We will see later another, more general way of getting the energy flow.)

24-5 Observing guided waves

Energy can be coupled into a waveguide by some kind of an “antenna.” For example, a little vertical wire or “stub” will do. The presence of the guided waves can be observed by picking up some of the electromagnetic energy with a little receiving “antenna,” which again can be a little stub of wire or a small loop.

Fig. 24-8. A waveguide with a driving stub and a pickup probe.



In Fig. 24-8, we show a guide with some cutaways to show a driving stub and a pickup “probe”. The driving stub can be connected to a signal generator via a coaxial cable, and the pickup probe can be connected by a similar cable to a detector. It is usually convenient to insert the pickup probe via a long thin slot in the guide, as shown in Fig. 24-8. Then the probe can be moved back and forth along the guide to sample the fields at various positions.

If the signal generator is set at some frequency ω greater than the cutoff frequency ω_c , there will be waves propagated down the guide from the driving stub. These will be the only waves present if the guide is infinitely long, which can effectively be arranged by terminating the guide with a carefully designed absorber in such a way that there are no reflections from the far end. Then, since the detector measures the time average of the fields near the probe, it will pick up a signal which is independent of the position along the guide; its output will be proportional to the power being transmitted.

If now the far end of the guide is finished off in some way that produces a reflected wave—as an extreme example, if we closed it off with a metal plate—there will be a reflected wave in addition to the original forward wave. These two waves will interfere and produce a standing wave in the guide similar to the standing waves on a string which we discussed in Chapter 49 of Vol. I. Then, as the pickup probe is moved along the line, the detector reading will rise and fall periodically, showing a maximum in the fields at each loop of the standing wave and a minimum at each node. The distance between two successive nodes (or loops) is just $\lambda_g/2$. This gives a convenient way of measuring the guide wavelength. If the frequency is now moved closer to ω_c , the distances between nodes increase, showing that the guide wavelength increases as predicted by Eq. (24.19).

Suppose now the signal generator is set at a frequency just a little below ω_c . Then the detector output will decrease gradually as the pickup probe is moved down the guide. If the frequency is set somewhat lower, the field strength will fall rapidly, following the curve of Fig. 24-7, and showing that waves are not propagated.

24-6 Waveguide plumbing

An important practical use of waveguides is for the transmission of high-frequency power, as, for example, in coupling the high-frequency oscillator or output amplifier of a radar set to an antenna. In fact, the antenna itself usually consists of a parabolic reflector fed at its focus by a waveguide flared out at the end to make a “horn” that radiates the waves coming along the guide. Although high frequencies can be transmitted along a coaxial cable, a waveguide is better for transmitting large amounts of power. First, the maximum power that can be transmitted along a line is limited by the breakdown of the insulation (solid or gas) between the conductors. For a given amount of power, the field strengths in a guide are usually less than they are in a coaxial cable, so higher powers can be transmitted before breakdown occurs. Second, the power losses in the coaxial cable are usually greater than in a waveguide. In a coaxial cable there must be insulating material to support the central conductor, and there is an energy loss in this material—particularly at high frequencies. Also, the current densities on the central conductor are quite high, and since the losses go as the *square* of the current density, the lower currents that appear on the walls of the guide result in lower

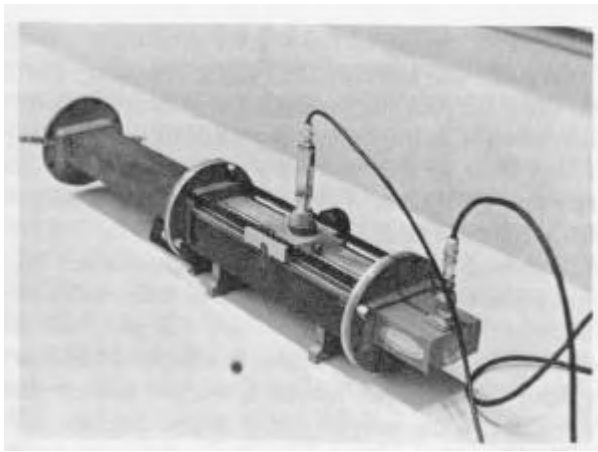


Fig. 24-9. Sections of waveguide connected with flanges.

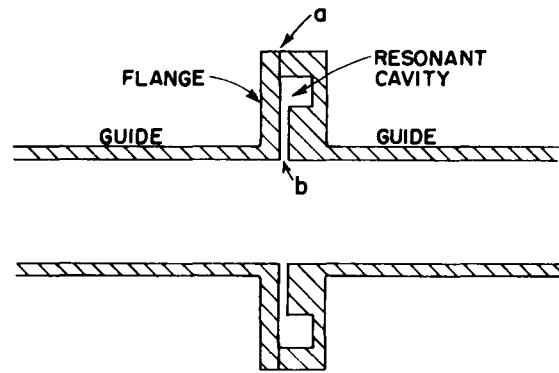


Fig. 24-10. A low-loss connection between two sections of waveguide.

energy losses. To keep these losses to a minimum, the inner surfaces of the guide are often plated with a material of high conductivity, such as silver.

The problem of connecting a "circuit" with waveguides is quite different from the corresponding circuit problem at low frequencies, and is usually called microwave "plumbing." Many special devices have been developed for the purpose. For instance, two sections of waveguide are usually connected together by means of flanges, as can be seen in Fig. 24-9. Such connections can, however, cause serious energy losses, because the surface currents must flow across the joint which may have a relatively high resistance. One way to avoid such losses is to make the flanges as shown in the cross section drawn in Fig. 24-10. A small space is left between the adjacent sections of the guide, and a groove is cut in the face of one of the flanges to make a small cavity of the type shown in Fig. 23-16(c). The dimensions are chosen so that this cavity is resonant at the frequency being used. This resonant cavity presents a high "impedance" to the currents, so relatively little current flows across the metallic joints (at *a* in Fig. 24-10). The high guide currents simply charge and discharge the "capacity" of the gap (at *b* in the figure), where there is little dissipation of energy.

Suppose you want to stop a waveguide in a way that won't result in reflected waves. Then you must put something at the end that imitates an infinite length of guide. You need a "termination" which acts for the guide like the characteristic impedance does for a transmission line—something that absorbs the arriving waves without making reflections. Then the guide will act as though it went on forever. Such terminations are made by putting inside the guide some wedges of resistance material carefully designed to absorb the wave energy while generating almost no reflected waves.

If you want to connect *three* things together—for instance, one source to two different antennas—then you can use a "T" like the one shown in Fig. 24-11. Power fed in at the center section of the "T" will be split and go out the two side arms (and there may also be some reflected waves). You can see qualitatively from the sketches in Fig. 24-12 that the fields would spread out when they get to the end of the input section and make electric fields that will start waves going out the two arms. Depending on whether electric fields in the guide are parallel or perpendicular to the "top" of the "T," the fields at the junction would be roughly as shown in (a) or (b) of Fig. 24-12.

Finally, we would like to describe a device called an "unidirectional coupler," which is very useful for telling what is going on after you have connected a complicated arrangement of waveguides. Suppose you want to know which way the waves are going in a particular section of guide—you might be wondering, for instance, whether or not there is a strong reflected wave. The unidirectional coupler takes out a small fraction of the power of a guide if there is a wave going one way, but none if the wave is going the other way. By connecting the output of the coupler to a detector, you can measure the "one-way" power in the guide

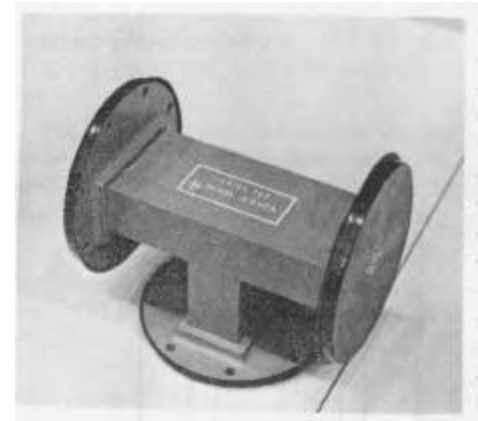


Fig. 24-11. A waveguide "T." (The flanges have plastic end caps to keep the inside clean while the "T" is not being used.)

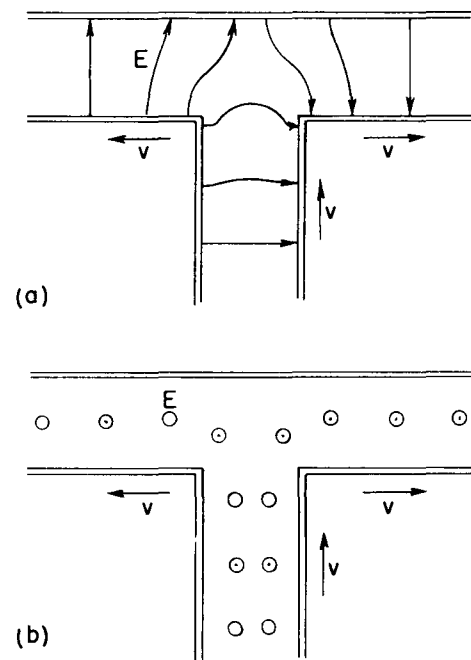


Fig. 24-12. The electric fields in a waveguide "T" for two possible field orientations.

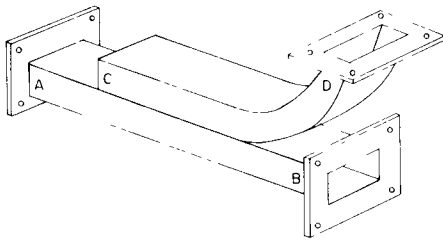


Fig 24-13. A unidirectional coupler.

Figure 24-13 is a drawing of a unidirectional coupler; a piece of waveguide AB has another piece of waveguide CD soldered to it along one face. The guide CD is curved away so that there is room for the connecting flanges. Before the guides are soldered together, two (or more) holes have been drilled in each guide (matching each other) so that some of the fields in the main guide AB can be coupled into the secondary guide CD . Each of the holes acts like a little antenna that produces a wave in the secondary guide. If there were only one hole, waves would be sent in both directions and would be the same no matter which way the wave was going in the primary guide. But when there are *two* holes with a separation space equal to one-quarter of the guide wavelength, they will make two sources 90° out of phase. Do you remember that we considered in Chapter 29 of Vol I the interference of the waves from two antennas spaced $\lambda/4$ apart and excited 90° out of phase in time? We found that the waves subtract in one direction and add in the opposite direction. The same thing will happen here. The wave produced in the guide CD will be going in the same direction as the wave in AB .

If the wave in the primary guide is travelling from A toward B , there will be a wave at the output D of the secondary guide. If the wave in the primary guide goes from B toward A , there will be a wave going toward the end C of the secondary guide. This end is equipped with a termination, so that this wave is absorbed and there is no wave at the output of the coupler.

24-7 Waveguide modes

The wave we have chosen to analyze is a special solution of the field equations. There are many more. Each solution is called a waveguide "mode." For example, our x -dependence of the field was just one-half a cycle of a sine wave. There is an equally good solution with a full cycle, then the variation of E_y with x is as shown in Fig 24-14. The k_x for such a mode is twice as large, so the cutoff frequency is much higher. Also, in the wave we studied E has only a y -component, but there are other modes with more complicated electric fields. If the electric field has components only in x and y —so that the total electric field is always at right angles to the z -direction—the mode is called a "transverse electric" (or TE) mode. The magnetic field of such modes will always have a z -component. It turns out that if E has a component in the z -direction (along the direction of propagation), then the magnetic field will always have only transverse components. So such fields are called transverse magnetic (TM) modes. For a rectangular guide, all the other modes have a higher cutoff frequency than the simple TE mode we have described. It is, therefore, possible—and usual—to use a guide with a frequency just above the cutoff for this lowest mode but below the cutoff frequency for all the others, so that just the one mode is propagated. Otherwise, the behavior gets complicated and difficult to control.

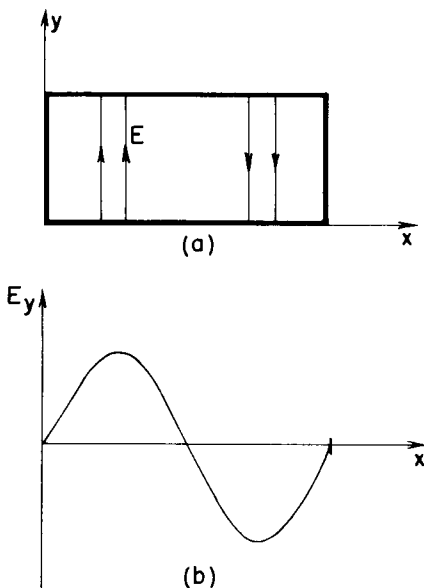


Fig. 24-14. Another possible variation of E_y with x .

24-8 Another way of looking at the guided waves

We want now to show you another way of understanding why a waveguide attenuates the fields rapidly for frequencies below the cutoff frequency ω_c . Then you will have a more "physical" idea of why the behavior changes so drastically between low and high frequencies. We can do this for the rectangular guide by analyzing the fields in terms of reflections—or images—in the walls of the guide. The approach only works for rectangular guides, however, that's why we started with the more mathematical analysis which works, in principle, for guides of any shape.

For the mode we have described, the vertical dimension (in y) had no effect, so we can ignore the top and bottom of the guide and imagine that the guide is extended indefinitely in the vertical direction. We imagine then that the guide just consists of two vertical plates with the separation a .

Let's say that the source of the fields is a vertical wire placed in the middle of the guide, with the wire carrying a current that oscillates at the frequency ω . In the absence of the guide walls such a wire would radiate cylindrical waves

Now we consider that the guide walls are perfect conductors. Then, just as in electrostatics, the conditions at the surface will be correct if we add to the field of the wire the field of one or more suitable image wires. The image idea works just as well for electrodynamics as it does for electrostatics, provided, of course, that we also include the retardations. We know that is true because we have often seen a mirror producing an image of a light source. And a mirror is just a "perfect" conductor for electromagnetic waves with optical frequencies.

Now let's take a horizontal cross section, as shown in Fig. 24-15, where W_1 and W_2 are the two guide walls and S_0 is the source wire. We call the direction of the current in the wire positive. Now if there were only one wall, say W_1 , we could remove it if we placed an image source (with opposite polarity) at the position marked S_1 . But with both walls in place there will also be an image of S_0 in the wall W_2 , which we show as the image S_2 . This source, too, will have an image in W_1 , which we call S_3 . Now both S_1 and S_3 will have images in W_2 at the positions marked S_4 and S_6 , and so on. For our two plane conductors with the source halfway between, the fields are the same as those produced by an infinite line of sources, all separated by the distance a . (It is, in fact just what you would see if you looked at a wire placed halfway between two parallel mirrors.) For the fields to be zero at the walls, the polarity of the currents in the images must alternate from one image to the next. In other words, they oscillate 180° out of phase. The waveguide field is, then, just the superposition of the fields of such an infinite set of line sources.

We know that if we are close to the sources, the field is very much like the static fields. We considered in Section 7-5 the static field of a grid of line sources and found that it is like the field of a charged plate except for terms that decrease exponentially with the distance from the grid. Here the average source strength is zero, because the sign alternates from one source to the next. Any fields which exist should fall off exponentially with distance. Close to the source, we see the field mainly of the nearest source; at large distances, many sources contribute and their average effect is zero. So now we see why the waveguide below cutoff frequency gives an exponentially decreasing field. At low frequencies, in particular, the static approximation is good, and it predicts a rapid attenuation of the fields with distance.

Now we are faced with the opposite question: Why are waves propagated at all? That is the mysterious part! The reason is that at high frequencies the retardation of the fields can introduce additional changes in phase which can cause the fields of the out-of-phase sources to add instead of cancelling. In fact, in Chapter 29 of Vol. I we have already studied, just for this problem, the fields generated by an array of antennas or by an optical grating. There we found that when several radio antennas are suitably arranged, they can give an interference pattern that has a strong signal in some direction but no signal in another.

Suppose we go back to Fig. 24-15 and look at the fields which arrive at a large distance from the array of image sources. The fields will be strong only in certain directions which depend on the frequency—only in those directions for which the fields from all the sources add in phase. At a reasonable distance from the sources the field propagates in these special directions as plane waves. We have sketched such a wave in Fig. 24-16, where the solid lines represent the wave crests and the dashed lines represent the troughs. The wave direction will be the one for which the difference in the retardation for two neighboring sources to the crest of a wave corresponds to one-half a period of oscillation. In other words, the difference between r_2 and r_0 in the figure is one-half of the free-space wavelength:

$$r_2 - r_0 = \frac{\lambda_0}{2}.$$

The angle θ is then given by

$$\sin \theta = \frac{\lambda_0}{2a}. \quad (24.33)$$

There is, of course, another set of waves travelling downward at the symmetric angle with respect to the array of sources. The complete waveguide field (not too

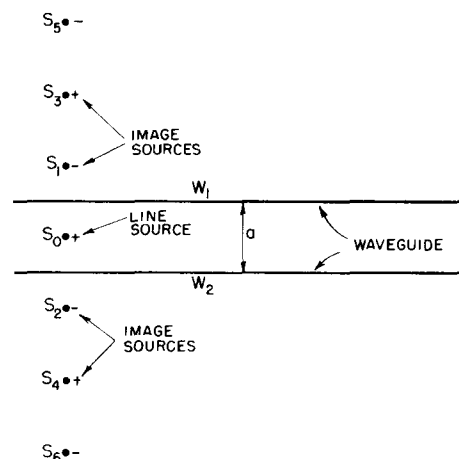


Fig. 24-15. The line source S_0 between the conducting plane walls W_1 and W_2 . The walls can be replaced by the infinite sequence of image sources.

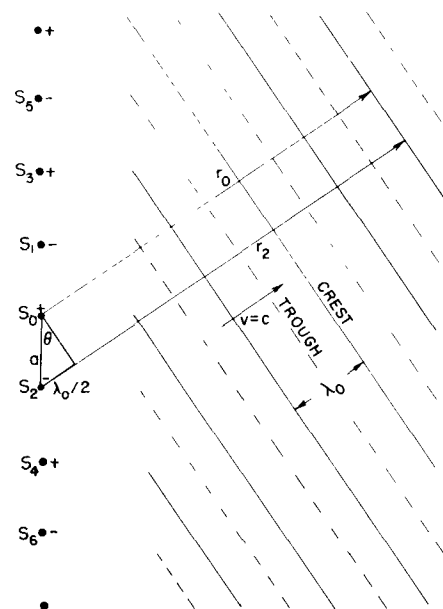


Fig. 24-16. One set of coherent waves from an array of line sources.

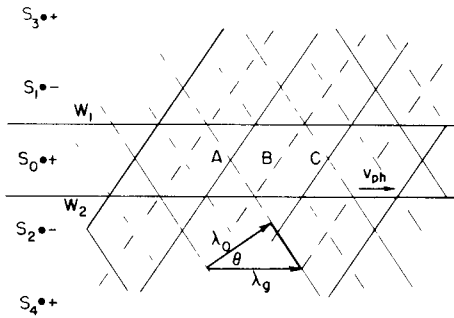


Fig. 24-17. The waveguide field can be viewed as the superposition of two trains of plane waves.

close to the source) is the superposition of these two sets of waves, as shown in Fig. 24-17. The actual fields are really like this, of course, only between the two walls of the waveguide.

At points like *A* and *C*, the crests of the two wave patterns coincide, and the field will have a maximum, at points like *B*, both waves have their peak negative value, and the field has its minimum (largest negative) value. As time goes on the field in the guide appears to be travelling along the guide with a wavelength λ_g , which is the distance from *A* to *C*. That distance is related to θ by

$$\cos \theta = \frac{\lambda_0}{\lambda_g}. \quad (24.34)$$

Using Eq. (24.33) for θ , we get that

$$\lambda_g = \frac{\lambda_0}{\cos \theta} = \frac{\lambda_0}{\sqrt{1 - (\lambda_0/2a)^2}}, \quad (24.35)$$

which is just what we found in Eq. (24.19).

Now we see why there is only wave propagation above the cutoff frequency ω_0 . If the free-space wavelength is longer than $2a$, there is no angle where the waves shown in Fig. 24-16 can appear. The necessary constructive interference appears suddenly when λ_0 drops below $2a$, or when ω goes above $\omega_0 = \pi c/a$.

If the frequency is high enough, there can be two or more possible directions in which the waves will appear. For our case, this will happen if $\lambda_0 < \frac{2}{3}a$. In general, however, it could also happen when $\lambda_0 < a$. These additional waves correspond to the higher guide modes we have mentioned.

It has also been made evident by our analysis why the phase velocity of the guided waves is greater than c and why this velocity depends on ω . As ω is changed, the angle of the free waves of Fig. 24-16 changes, and therefore so does the velocity along the guide.

Although we have described the guided wave as the superposition of the fields of an infinite array of line sources, you can see that we would arrive at the same result if we imagined two sets of free-space waves being continually reflected back and forth between two perfect mirrors—remembering that a reflection means a reversal of phase. These sets of reflecting waves would all cancel each other unless they were going at just the angle θ given in Eq. (24.33). There are many ways of looking at the same thing.

Electrodynamics in Relativistic Notation

25-1 Four-vectors

We now discuss the application of the special theory of relativity to electrodynamics. Since we have already studied the special theory of relativity in Chapters 15 through 17 of Vol. I, we will just review quickly the basic ideas.

It is found experimentally that the laws of physics are unchanged if we move with uniform velocity. You can't tell if you are inside a spaceship moving with uniform velocity in a straight line, unless you look outside the spaceship, or at least make an observation having to do with the world outside. Any true law of physics we write down must be arranged so that this fact of nature is built in.

The relationship between the space and time of two systems of coordinates, one, S' , in uniform motion in the x -direction with speed v relative to the other, S , is given by the *Lorentz transformation*:

$$\begin{aligned} t' &= \frac{t - vx}{\sqrt{1 - v^2}}, & y' &= y, \\ x' &= \frac{x - vt}{\sqrt{1 - v^2}}, & z' &= z. \end{aligned} \quad (25.1)$$

The laws of physics must be such that after a Lorentz transformation, the new form of the laws looks just like the old form. This is just like the principle that the laws of physics don't depend on the *orientation* of our coordinate system. In Chapter 11 of Vol. I, we saw that the way to describe mathematically the invariance of physics with respect to rotations was to write our equations in terms of *vectors*.

For example, if we have two vectors

$$\mathbf{A} = (A_x, A_y, A_z) \quad \text{and} \quad \mathbf{B} = (B_x, B_y, B_z),$$

we found that the combination

$$\mathbf{A} \cdot \mathbf{B} = A_x B_x + A_y B_y + A_z B_z$$

was not changed if we transformed to a rotated coordinate system. So we know that if we have a scalar product like $\mathbf{A} \cdot \mathbf{B}$ on both sides of an equation, the equation will have exactly the same form in all rotated coordinate systems. We also discovered an operator (see Chapter 2),

$$\nabla = \left(\frac{\partial}{\partial x}, \frac{\partial}{\partial y}, \frac{\partial}{\partial z} \right),$$

which, when applied to a scalar function, gave three quantities which transform just like a vector. With this operator we defined the gradient, and in combination with other vectors, the divergence and the Laplacian. Finally we discovered that by taking sums of certain products of pairs of the components of two vectors we could get three new quantities which behaved like a new vector. We called it the *cross product* of two vectors. Using the cross product with our operator ∇ we then defined the curl of a vector.

Since we will be referring back to what we have done in vector analysis, we have put in Table 25-1 a summary of all the important vector operations in three dimensions that we have used in the past. The point is that it must be possible to write the equations of physics so that both sides transform the same way under

25-1 Four-vectors

25-2 The scalar product

25-3 The four-dimensional gradient

25-4 Electrodynamics in four-dimensional notation

25-5 The four-potential of a moving charge

25-6 The invariance of the equations of electrodynamics

In this chapter: $c = 1$

Review: Chapter 15, Vol. I, *The Special Theory of Relativity*
 Chapter 16, Vol. I, *Relativistic Energy and Momentum*
 Chapter 17, Vol. I, *Space-Time*
 Chapter 13, Vol. II, *Magnetostatics*

Table 25-1

The important quantities and operations of vector analysis in three dimensions

Definition of a vector	$\mathbf{A} = (A_x, A_y, A_z)$
Scalar product	$\mathbf{A} \cdot \mathbf{B}$
Differential vector operator	∇
Gradient	$\nabla\phi$
Divergence	$\nabla \cdot \mathbf{A}$
Laplacian	$\nabla \cdot \nabla = \nabla^2$
Cross product	$\mathbf{A} \times \mathbf{B}$
Curl	$\nabla \times \mathbf{A}$

rotations. If one side is a vector, the other side must also be a vector, and both sides will change together in exactly the same way if we rotate our coordinate system. Similarly, if one side is a scalar, the other side must also be a scalar, so that neither side changes when we rotate coordinates, and so on.

Now in the case of special relativity, time and space are inextricably mixed, and we must do the analogous things for four dimensions. We want our equations to remain the same not only for rotations, but also for *any* inertial frame. That means that our equations should be invariant under the Lorentz transformation of equations (25.1). The purpose of this chapter is to show you how that can be done. Before we get started, however, we want to do something that makes our work a lot easier (and saves some confusion). And that is to choose our units of length and time so that the speed of light c is equal to 1. You can think of it as taking our unit of time to be *the time that it takes light to go one meter* (which is about 3×10^{-9} sec). We can even call this time unit "one meter." Using this unit, all of our equations will show more clearly the space-time symmetry. Also, all the c 's will disappear from our relativistic equations. (If this bothers you, you can always put the c 's back into any equation by replacing every t by ct , or, in general, by sticking in a c wherever it is needed to make the dimensions of the equations come out right.) With this groundwork we are ready to begin. Our program is to do in the four dimensions of space-time all of the things we did with vectors for three dimensions. It is really quite a simple game, we just work by analogy. The only real complications is the notation (we've already used up the vector symbol for three dimensions) and one slight twist of signs.

First, by analogy with vectors in three dimensions, we define a *four-vector* as a set of the four quantities $a_t, a_x, a_y,$ and a_z , which transform like $t, x, y,$ and z when we change to a moving coordinate system. There are several different notations people use for a four-vector; we will write a_μ , by which we mean the group of four numbers (a_t, a_x, a_y, a_z)—in other words, the subscript μ can take on the four "values" t, x, y, z . It will also be convenient, at times, to indicate the three space components by a three-vector, like this: $a_\mu = (a_t, \mathbf{a})$.

We have already encountered one four-vector, which consists of the energy and momentum of a particle (Chapter 17, Vol. I). In our new notation we write

$$p_\mu = (E, \mathbf{p}), \quad (25.2)$$

which means that the four-vector p_μ is made up of the energy E and the three components of the three-vector \mathbf{p} of a particle.

It looks as though the game is really very simple—for each three-vector in physics all we have to do is find what the remaining component should be, and we have a four-vector. To see that this is not the case, consider the velocity vector with components

$$v_x = \frac{dx}{dt}, \quad v_y = \frac{dy}{dt}, \quad v_z = \frac{dz}{dt}.$$

The question is: What is the time component? Instinct should give the right answer. Since four-vectors are like t, x, y, z , we would guess that the time component is

$$v_t = \frac{dt}{dt} = 1.$$

This is wrong. The reason is that the t in each denominator is not an invariant when we make a Lorentz transformation. The numerators have the right behavior to make a four-vector, but the dt in the denominator spoils things; it is unsymmetric and is not the same in two different systems.

It turns out that the four "velocity" components which we have written down will become the components of a four-vector if we just divide by $\sqrt{1 - v^2}$. We can see that that is true because if we start with the momentum four-vector

$$p_\mu = (E, \mathbf{p}) = \left(\frac{m_0}{\sqrt{1 - v^2}}, \frac{m_0 \mathbf{v}}{\sqrt{1 - v^2}} \right), \quad (25.3)$$

and divide it by the rest mass m_0 , which is an invariant scalar in *four dimensions*, we have

$$\frac{p_\mu}{m_0} = \left(\frac{1}{\sqrt{1-v^2}}, \frac{\mathbf{v}}{\sqrt{1-v^2}} \right), \quad (24.4)$$

which must still be a four-vector. (Dividing by an *invariant scalar* doesn't change the transformation properties) So we can *define* the “velocity four-vector” u_μ by

$$\begin{aligned} u_t &= \frac{1}{\sqrt{1-v^2}}, & u_y &= \frac{v_y}{\sqrt{1-v^2}}, \\ u_x &= \frac{v_x}{\sqrt{1-v^2}}, & u_z &= \frac{v_z}{\sqrt{1-v^2}}, \end{aligned} \quad (25.5)$$

The four-velocity is a useful quantity; we can, for instance, write

$$p_\mu = m_0 u_\mu. \quad (25.6)$$

This is the typical sort of form an equation which is relativistically correct must have; each side is a four-vector. (The right-hand side is an invariant times a four-vector, which is still a four-vector.)

25-2 The scalar product

It is an accident of life, if you wish, that under coordinate rotations the distance of a point from the origin does not change. This means mathematically that $r^2 = x^2 + y^2 + z^2$ is an invariant. In other words, after a rotation $r'^2 = r^2$, or

$$x'^2 + y'^2 + z'^2 = x^2 + y^2 + z^2.$$

Now the question is: Is there a similar quantity which is invariant under the Lorentz transformation? There is. From Eq. (25.1) you can see that

$$t'^2 - x'^2 = t^2 - x^2$$

That is pretty nice, except that it depends on a particular choice of the x -direction. We can fix that up by subtracting y^2 and z^2 . Then any Lorentz transformation *plus* a rotation will leave the quantity unchanged. So the quantity which is analogous to r^2 for four dimensions, in three dimensions is

$$t^2 - x^2 - y^2 - z^2.$$

It is an invariant under what is called the “complete Lorentz group”—which means for transformation of both translations at constant velocity *and* rotations.

Now since this invariance is an algebraic matter depending only on the transformation rules of Eq. (25.1)—plus rotations—it is true for any four-vector (by definition they all transform the same). So for a four-vector a_μ we have that

$$a_t'^2 - a_x'^2 - a_y'^2 - a_z'^2 = a_t^2 - a_x^2 - a_y^2 - a_z^2.$$

We will call this quantity the square of “the length” of the four-vector a_μ . (Sometimes people change the sign of all the terms and call the length $a_x^2 + a_y^2 + a_z^2 - a_t^2$, so you'll have to watch out.)

Now if we have *two* vectors a_μ and b_μ , their corresponding components transform in the same way, so the combination

$$a_t b_t - a_x b_x - a_y b_y - a_z b_z$$

is also an invariant (scalar) quantity. (We have in fact already proved this in Chapter 17 of Vol. I.) Clearly this expression is quite analogous to the dot product for vectors. We will, in fact, call it the *dot product* or *scalar product* of two four-vectors. It would seem logical to write it as $a_\mu \cdot b_\mu$, so it would *look* like a dot product. But, unhappily, it's not done that way; it is usually written without the dot.

So we will follow the convention and write the dot product simply as $a_\mu b_\mu$. So, by definition,

$$a_\mu b_\mu = a_t b_t - a_x b_x - a_y b_y - a_z b_z \quad (25.7)$$

Whenever you see two identical subscripts together (we will occasionally have to use ν or some other letter instead of μ) it means that you are to take the four products and sum, remembering the minus sign for the products of the space components. With this convention the invariance of the scalar product under a Lorentz transformation can be written as

$$a'_\mu b'_\mu = a_\mu b_\mu.$$

Since the last three terms in (25.7) are just the scalar dot product in three dimensions, it is often more convenient to write

$$a_\mu b_\mu = a_t b_t - \mathbf{a} \cdot \mathbf{b}.$$

It is also obvious that the four-dimensional length we described above can be written as $a_\mu a_\mu$:

$$a_\mu a_\mu = a_t^2 - a_x^2 - a_y^2 - a_z^2 = a_t^2 - \mathbf{a} \cdot \mathbf{a}. \quad (25.8)$$

It will also be convenient to sometimes write this quantity as a_μ^2 :

$$a_\mu^2 \equiv a_\mu a_\mu.$$

We will now give you an illustration of the usefulness of four-vector dot products. Antiprotons (\bar{P}) are produced in large accelerators by the reaction

$$P + P \rightarrow P + P + P + \bar{P}.$$

That is, an energetic proton collides with a proton at rest (for example, in a hydrogen target placed in the beam), and if the incident proton has enough energy, a proton-antiproton pair may be produced, in addition to the two original protons.* The question is: How much energy must be given to the incident proton to make this reaction energetically possible?

The easiest way to get the answer is to consider what the reaction looks like in the center-of-mass (CM) system (see Fig. 25-1). We'll call the incident proton a and its four-momentum p_μ^a . Similarly, we'll call the target proton b and its four-

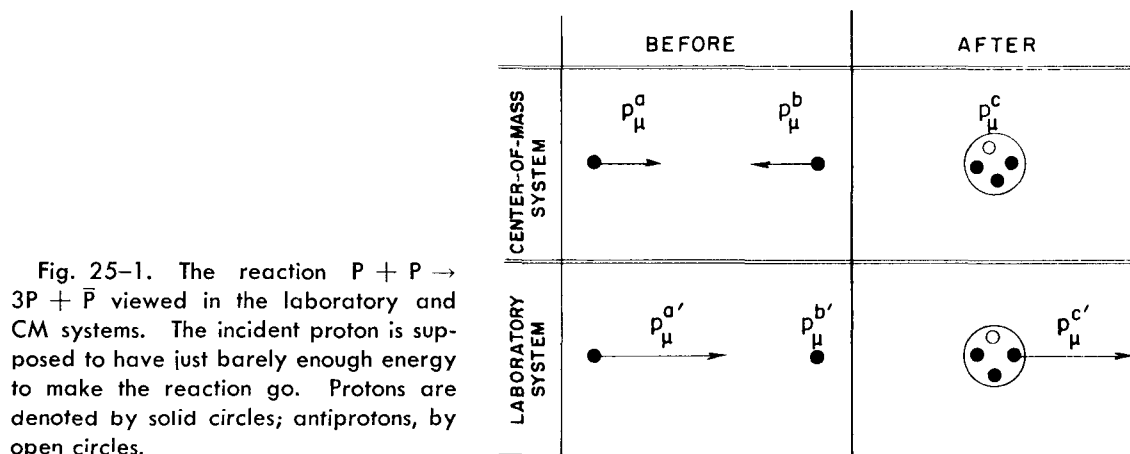


Fig. 25-1. The reaction $P + P \rightarrow 3P + \bar{P}$ viewed in the laboratory and CM systems. The incident proton is supposed to have just barely enough energy to make the reaction go. Protons are denoted by solid circles; antiprotons, by open circles.

* You may well ask: Why not consider the reactions

$$P + P \rightarrow P + P + \bar{P},$$

or even

$$P + P \rightarrow P + \bar{P}$$

which clearly require less energy? The answer is that a principle called *conservation of baryons* tells us the quantity "number of protons minus number of antiprotons" cannot change. This quantity is 2 on the left side of our reaction. Therefore, if we want an antiproton on the right side, we must have also *three* protons (or other baryons).

momentum p_μ^b . If the incident proton has *just barely* enough energy to make the reaction go, the final state—the situation after the collision—will consist of a glob containing three protons and an antiproton at rest in the CM system. If the incident energy were slightly higher, the final state particles would have some kinetic energy and be moving apart; if the incident energy were slightly lower, there would not be enough energy to make the four particles.

If we call p_μ^c the total four-momentum of the whole glob in the final state, conservation of energy and momentum tells us that

$$\mathbf{p}^a + \mathbf{p}^b = \mathbf{p}^c,$$

and

$$E^a + E^b = E^c.$$

Combining these two equations, we can write that

$$p_\mu^a + p_\mu^b = p_\mu^c. \quad (25.9)$$

Now the important thing is that this is an equation among four-vectors, and is, therefore, true in any inertial frame. We can use this fact to simplify our calculations. We start by taking the “length” of each side of Eq. (25.9); they are, of course, also equal. We get

$$(p_\mu^a + p_\mu^b)(p_\mu^a + p_\mu^b) = p_\mu^c p_\mu^c. \quad (25.10)$$

Since $p_\mu^c p_\mu^c$ is invariant, we can evaluate it in any coordinate system. In the CM system, the time component of p_μ^c is the rest energy of four protons, namely $4M$, and the space part \mathbf{p} is zero; so $p_\mu^c = (4M, \mathbf{0})$. We have used the fact that the rest mass of an antiproton equals the rest mass of a proton, and we have called this common mass M .

Thus, Eq. (25.10) becomes

$$p_\mu^a p_\mu^a + 2p_\mu^a p_\mu^b + p_\mu^b p_\mu^b = 16M^2. \quad (25.11)$$

Now $p_\mu^a p_\mu^a$ and $p_\mu^b p_\mu^b$ are very easy, since the “length” of the momentum four-vector of any particle is just the mass of the particle squared:

$$p_\mu p_\mu = E^2 - \mathbf{p}^2 = M^2.$$

This can be shown by direct calculation or, more cleverly, by noting that for a particle *at rest* $p_\mu = (M, \mathbf{0})$, so $p_\mu p_\mu = M^2$. But since it is an invariant, it is equal to M^2 in *any* frame. Using these results in Eq. (25.11), we have

$$2p_\mu^a p_\mu^b = 14M^2$$

or

$$p_\mu^a p_\mu^b = 7M^2. \quad (25.12)$$

Now we can also evaluate $p_\mu^a p_\mu^b$ in the laboratory system. The four-vector p_μ^a can be written (E^a, \mathbf{p}^a) , while $p_\mu^b = (M, \mathbf{0})$, since it describes a proton at rest. Thus, $p_\mu^a p_\mu^b$ must also be equal to ME^a , and since we know the scalar product is an invariant this must be numerically the same as what we found in (25.12). So we have that

$$E^a = 7M,$$

which is the result we were after. The *total* energy of the initial proton must be at least $7M$ (about 6.6 Gev since $M = 938$ Mev) or, subtracting the rest mass M , the *kinetic* energy must be at least $6M$ (about 5.6 Gev). The Bevatron accelerator at Berkeley was designed to give about 6.2 Gev of kinetic energy to the protons it accelerates, in order to be able to make antiprotons.

Since scalar products are invariant, they are always interesting to evaluate. What about the “length” of the four-velocity $u_\mu u_\mu$?

$$u_\mu u_\mu = u_t^2 - \mathbf{u}^2 = \frac{1}{1 - v^2} - \frac{v^2}{1 - v^2} = 1.$$

Thus, u_μ is the *unit four-vector*

25-3 The four-dimensional gradient

The next thing that we have to discuss is the four-dimensional analog of the gradient. We recall (Chapter 14, Vol. I) that the three differential operators $\partial/\partial x$, $\partial/\partial y$, $\partial/\partial z$ transform like a three-vector and are called the gradient. The same scheme ought to work in four dimensions; that is, we might guess that the four-dimensional gradient should be $(\partial/\partial t, \partial/\partial x, \partial/\partial y, \partial/\partial z)$. *This is wrong.*

To see the error, consider a scalar function ϕ which depends only on x and t . The change in ϕ , if we make a small change Δt in t while holding x constant, is

$$\Delta\phi = \frac{\partial\phi}{\partial t} \Delta t. \quad (25.13)$$

On the other hand, according to a moving observer,

$$\Delta\phi = \frac{\partial\phi}{\partial x'} \Delta x' + \frac{\partial\phi}{\partial t'} \Delta t'.$$

We can express $\Delta x'$ and $\Delta t'$ in terms of Δt by using Eq. (25.1). Remembering that we are holding x constant, so that $\Delta x = 0$, we write

$$\Delta x' = -\frac{v}{\sqrt{1-v^2}} \Delta t; \quad \Delta t' = \frac{\Delta t}{\sqrt{1-v^2}}.$$

Thus,

$$\begin{aligned} \Delta\phi &= \frac{\partial\phi}{\partial x'} \left(-\frac{v}{\sqrt{1-v^2}} \Delta t \right) + \frac{\partial\phi}{\partial t'} \left(\frac{\Delta t}{\sqrt{1-v^2}} \right) \\ &= \left(\frac{\partial\phi}{\partial t'} - v \frac{\partial\phi}{\partial x'} \right) \frac{\Delta t}{\sqrt{1-v^2}}. \end{aligned}$$

Comparing this result with Eq. (25.13), we learn that

$$\frac{\partial\phi}{\partial t} = \frac{1}{\sqrt{1-v^2}} \left(\frac{\partial\phi}{\partial t'} - v \frac{\partial\phi}{\partial x'} \right) \quad (25.14)$$

A similar calculation gives

$$\frac{\partial\phi}{\partial x} = \frac{1}{\sqrt{1-v^2}} \left(\frac{\partial\phi}{\partial x'} - v \frac{\partial\phi}{\partial t'} \right). \quad (25.15)$$

Now we can see that the gradient is rather strange. The formulas for x and t in terms of x' and t' [obtained by solving Eq. (25.1)] are:

$$t = \frac{t' + vx'}{\sqrt{1-v^2}}, \quad x = \frac{x' + vt'}{\sqrt{1-v^2}}.$$

This is the way a four-vector *must* transform. But Eqs. (25.14) and (25.15) have a couple of signs wrong!

The answer is that instead of the *incorrect* $(\partial/\partial t, \nabla)$, we must *define* the *four-dimensional gradient* operator, which we will call ∇_μ , by

$$\nabla_\mu = \left(\frac{\partial}{\partial t}, -\nabla \right) = \left(\frac{\partial}{\partial t}, -\frac{\partial}{\partial x}, -\frac{\partial}{\partial y}, -\frac{\partial}{\partial z} \right). \quad (25.16)$$

With this definition, the sign difficulties encountered above go away, and ∇_μ behaves as a four-vector should. (It's rather awkward to have those minus signs, but that's the way the world is.) Of course, what it means to say that ∇_μ "behaves like a four-vector" is simply that the four-gradient of a scalar is a four-vector. If ϕ is a true scalar invariant field (Lorentz invariant) then $\nabla_\mu \phi$ is a four-vector field.

All right, now that we have vectors, gradients, and dot products, the next thing is to look for an invariant which is analogous to the divergence of three-dimensional vector analysis. Clearly, the analog is to form the expression $\nabla_\mu b_\mu$, where b_μ is a four-vector field whose components are functions of space and time.

We define the *divergence* of the four-vector $b_\mu = (b_t, \mathbf{b})$ as the dot product of ∇_μ and b_μ :

$$\begin{aligned}\nabla_\mu b_\mu &= \frac{\partial}{\partial t} b_t - \left(-\frac{\partial}{\partial x}\right) b_x - \left(-\frac{\partial}{\partial y}\right) b_y - \left(-\frac{\partial}{\partial z}\right) b_z \\ &= \frac{\partial}{\partial t} b_t + \nabla \cdot \mathbf{b},\end{aligned}\quad (25.17)$$

where $\nabla \cdot \mathbf{b}$ is the ordinary three-divergence of the three-vector \mathbf{b} . Note that one has to be careful with the signs. Some of the minus signs come from the definition of the scalar product, Eq. (25.7); the others are required because the space components of ∇_μ are $-\partial/\partial x$, etc., as in Eq. (25.16). The divergence as defined by (25.17) is an invariant and gives the same answer in all coordinate systems which differ by a Lorentz transformation.

Let's look at a physical example in which the four-divergence shows up. We can use it to solve the problem of the fields around a moving wire. We have already seen (Section 13-7) that the electric charge density ρ and the current density \mathbf{j} form a four-vector $j_\mu = (\rho, \mathbf{j})$. If an uncharged wire carries the current j_x , then in a frame moving past it with velocity v (along x), the wire will have the charge and current density [obtained from the Lorentz transformation Eqs. (25.1)] as follows:

$$\rho' = \frac{-vj_x}{\sqrt{1-v^2}}, \quad j'_x = \frac{j_x}{\sqrt{1-v^2}}.$$

These are just what we found in Chapter 13. We can then use these sources in Maxwell's equation in *the moving system* to find the fields.

The charge conservation law, Section 13-2, also takes on a simple form in the four-vector notation. Consider the four-divergence of j_μ :

$$\nabla_\mu j_\mu = \frac{\partial \rho}{\partial t} + \nabla \cdot \mathbf{j}. \quad (25.18)$$

The law of the conservation of charge says that the outflow of current per unit volume must equal the negative rate of increase of charge density. In other words, that

$$\nabla \cdot \mathbf{j} = -\frac{\partial \rho}{\partial t}.$$

Putting this into Eq. (25.18), the law of conservation of charge takes on the simple form

$$\nabla_\mu j_\mu = 0. \quad (25.19)$$

Since $\nabla_\mu j_\mu$ is an invariant scalar, if it is zero in one frame it is zero in all frames. We have the result that if charge is conserved in one coordinate system, it is conserved in all coordinate systems moving with uniform velocity.

As our last example we want to consider the scalar product of the gradient operator ∇_μ with itself. In three dimensions, such a product gives the Laplacian

$$\nabla^2 = \nabla \cdot \nabla = \frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} + \frac{\partial^2}{\partial z^2}.$$

What do we get in four dimensions? That's easy. Following our rules for dot products and gradients, we get

$$\begin{aligned}\nabla_\mu \nabla_\mu &= \frac{\partial}{\partial t} \frac{\partial}{\partial t} - \left(-\frac{\partial}{\partial x}\right) \left(-\frac{\partial}{\partial x}\right) - \left(-\frac{\partial}{\partial y}\right) \left(-\frac{\partial}{\partial y}\right) - \left(-\frac{\partial}{\partial z}\right) \left(-\frac{\partial}{\partial z}\right) \\ &= \frac{\partial^2}{\partial t^2} - \nabla^2.\end{aligned}$$

This operator, which is the analog of the three-dimensional Laplacian, is called

the *D'Alembertian* and has a special notation:

$$\square^2 = \nabla_\mu \nabla_\mu = \frac{\partial^2}{\partial t^2} - \nabla^2. \quad (25.20)$$

From its definition it is an invariant scalar operator; if it operates on a four-vector field, it produces a new four-vector field. (Some people define the D'Alembertian with the opposite sign to Eq. (25.20), so you will have to be careful when reading the literature.)

We have now found four-dimensional equivalents of most of the three-dimensional quantities we had listed in Table 25-1. (We do not yet have the equivalents of the cross product and the curl operation; we won't get to them until the next chapter.) It may help you remember how they go if we put all the important definitions and results together in one place, so we have made such a summary in Table 25-2.

Table 25-2

The important quantities of vector analysis in three and four dimensions.

	Three dimensions	Four dimensions
Vector	$\mathbf{A} = (A_x, A_y, A_z)$	$a_\mu = (a_t, a_x, a_y, a_z) = (a_t, \mathbf{a})$
Scalar product	$\mathbf{A} \cdot \mathbf{B} = A_x B_x + A_y B_y + A_z B_z$	$a_\mu b_\mu = a_t b_t - a_x b_x - a_y b_y - a_z b_z = a_t b_t - \mathbf{a} \cdot \mathbf{b}$
Vector operator	$\nabla = (\partial/\partial x, \partial/\partial y, \partial/\partial z)$	$\nabla_\mu = (\partial/\partial t, -\partial/\partial x, -\partial/\partial y, -\partial/\partial z) = (\partial/\partial t, -\nabla)$
Gradient	$\nabla\psi = \left(\frac{\partial\psi}{\partial x}, \frac{\partial\psi}{\partial y}, \frac{\partial\psi}{\partial z}\right)$	$\nabla_\mu\varphi = \left(\frac{\partial\varphi}{\partial t}, -\frac{\partial\varphi}{\partial x}, -\frac{\partial\varphi}{\partial y}, -\frac{\partial\varphi}{\partial z}\right) = \left(\frac{\partial\varphi}{\partial t}, \nabla\varphi\right)$
Divergence	$\nabla \cdot \mathbf{A} = \frac{\partial A_x}{\partial x} + \frac{\partial A_y}{\partial y} + \frac{\partial A_z}{\partial z}$	$\nabla_\mu a_\mu = \frac{\partial a_t}{\partial t} + \frac{\partial a_x}{\partial x} + \frac{\partial a_y}{\partial y} + \frac{\partial a_z}{\partial z} = \frac{\partial a_t}{\partial t} + \nabla \cdot \mathbf{a}$
Laplacian and D'Alembertian	$\nabla \cdot \nabla = \frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} + \frac{\partial^2}{\partial z^2}$	$\nabla_\mu \nabla_\mu = \frac{\partial^2}{\partial t^2} - \frac{\partial^2}{\partial x^2} - \frac{\partial^2}{\partial y^2} - \frac{\partial^2}{\partial z^2} = \frac{\partial^2}{\partial t^2} - \nabla^2 = \square^2$

25-4 Electrodynamics in four-dimensional notation

We have already encountered the D'Alembertian operator, without giving it that name, in Section 18-6; the differential equations we found there for the potentials can be written in the new notations as:

$$\square^2 \phi = \frac{\rho}{\epsilon_0}, \quad \square^2 \mathbf{A} = \frac{\mathbf{j}}{\epsilon_0}. \quad (25.21)$$

The four quantities on the right-hand side of the two equations in (25.21) are ρ, j_x, j_y, j_z , divided by ϵ_0 , which is a universal constant which will be the same in all coordinate systems if the same unit of charge is used in all frames. So the four quantities $\rho/\epsilon_0, j_x/\epsilon_0, j_y/\epsilon_0, j_z/\epsilon_0$ also transform as a four-vector. We can write them as j_μ/ϵ_0 . The D'Alembertian doesn't change when the coordinate system is changed, so the quantities ϕ, A_x, A_y, A_z *must also transform* like a four-vector—which means that they *are* the components of a four-vector. In short,

$$A_\mu = (\phi, \mathbf{A})$$

is a four-vector. What we call the scalar and vector potentials are really different aspects of the same physical thing. They belong together. And if they are kept together the relativistic invariance of the world is obvious. We call A_μ the *four-potential*.

In the four-vector notation Eqs. (25.21) become simply

$$\square^2 A_\mu = \frac{j_\mu}{\epsilon_0}, \quad (25.22)$$

The physics of this equation is just the same as Maxwell's equations. But there is some pleasure in being able to rewrite them in an elegant form. The pretty form is also meaningful, it shows directly the invariance of electrodynamics under the Lorentz transformation.

Remember that Eqs (25.21) could be deduced from Maxwell's equations only if we imposed the gauge condition

$$\frac{\partial \phi}{\partial t} + \nabla \cdot \mathbf{A} = 0, \quad (25.23)$$

which just says $\nabla_\mu A_\mu = 0$; the gauge condition says that the divergence of the four-vector A_μ is zero. This condition is called the *Lorentz condition*. It is very convenient because it is an invariant condition and therefore Maxwell's equations stay in the form of Eq. (25.22) for all frames.

25-5 The four-potential of a moving charge

Although it is implicit in what we have already said, let us write down the transformation laws which give ϕ and \mathbf{A} in a moving system in terms of ϕ and \mathbf{A} in a stationary system. Since $A_\mu = (\phi, \mathbf{A})$ is a four-vector, the equations must look just like Eqs. (25.1), except that t is replaced by ϕ , and \mathbf{x} is replaced by \mathbf{A} . Thus,

$$\begin{aligned} \phi' &= \frac{\phi - vA_x}{\sqrt{1 - v^2}}, & A'_y &= A_y, \\ A'_x &= \frac{A_x - v\phi}{\sqrt{1 - v^2}}, & A'_z &= A_z. \end{aligned} \quad (25.24)$$

This assumes that the primed coordinate system is moving with speed v in the positive x -direction, as measured in the unprimed coordinate system.

We will consider one example of the usefulness of the idea of the four-potential. What are the vector and scalar potentials of a charge q moving with speed v along the x -axis? The problem is easy in a coordinate system moving with the charge, since in this system the charge is standing still. Let's say that the charge is at the origin of the S' -frame, as shown in Fig. 25-2. The scalar potential in the moving system is then given by

$$\phi' = \frac{q}{4\pi\epsilon_0 r'}, \quad (25.25)$$

r' being the distance from q to the field point, as measured in the moving system. The vector potential \mathbf{A}' is, of course, zero.

Now it is straightforward to find ϕ and \mathbf{A} , the potentials as measured in the stationary coordinates. The inverse relations to Eqs. (25.24) are

$$\begin{aligned} \phi &= \frac{\phi' + vA'_x}{\sqrt{1 - v^2}}, & A_y &= A'_y, \\ A_x &= \frac{A'_x + v\phi'}{\sqrt{1 - v^2}}, & A_z &= A'_z. \end{aligned} \quad (25.26)$$

Using the ϕ' given by Eq. (25.25), and $\mathbf{A}' = 0$, we get

$$\begin{aligned} \phi &= \frac{q}{4\pi\epsilon_0} \frac{1}{r' \sqrt{1 - v^2}} \\ &= \frac{q}{4\pi\epsilon_0} \frac{1}{\sqrt{1 - v^2} \sqrt{x'^2 + y'^2 + z'^2}}. \end{aligned}$$

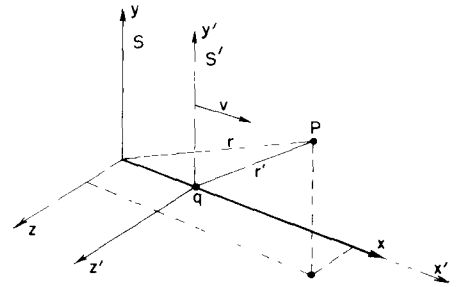


Fig. 25-2. The frame S' moves with velocity v (in the x -direction) with respect to S . A charge at rest at the origin of S' is at $x = vt$ in S . The potentials at P can be computed in either frame.

This gives us the scalar potential ϕ we would see in S , but, unfortunately, expressed in terms of the S' coordinates. We can get things in terms of t, x, y, z by substituting for t', x', y' , and z' , using (25.1). We get

$$\phi = \frac{q}{4\pi\epsilon_0} \frac{1}{\sqrt{1-v^2}} \frac{1}{\sqrt{[(x-vt)/\sqrt{1-v^2}]^2 + y^2 + z^2}}. \quad (25.27)$$

Following the same procedure for the components of A , you can show that

$$A = v\phi. \quad (25.28)$$

These are the same formulas we derived by a different method in Chapter 21.

25-6 The invariance of the equations of electrodynamics

We have found that the potentials ϕ and A taken together form a four-vector which we call A_μ , and that the wave equations—the full equations which determine the A_μ in terms of the j_μ —can be written as in Eq. (25.22). This equation, together with the conservation of charge, Eq. (25.19), gives us the fundamental law of the electromagnetic field:

$$\square^2 A_\mu = \frac{1}{\epsilon_0} j_\mu, \quad \nabla_\mu j_\mu = 0. \quad (25.29)$$

There, in one tiny space on the page, are all of the Maxwell equations—beautiful and simple. Did we learn anything from writing the equations this way, besides that they are beautiful and simple? In the first place, is it anything different from what we had before when we wrote everything out in all the various components? Can we from this equation deduce something that could not be deduced from the wave equations for the potentials in terms of the charges and currents? The answer is definitely no. The only thing we have been doing is changing the names of things—using a new notation. We have written a square symbol to represent the derivatives, but it still means nothing more nor less than the second derivative with respect to t , minus the second derivative with respect to x , minus the second derivative with respect to y , minus the second derivative with respect to z . And the μ means that we have four equations, one each for $\mu = t, x, y$, or z . What then is the significance of the fact that the equations can be written in this simple form? From the point of view of deducing anything directly, it doesn't mean anything. Perhaps, though, the simplicity of the equations means that nature also has a certain simplicity.

Let us show you something interesting that we have recently discovered: *All of the laws of physics can be contained in one equation.* That equation is

$$U = 0. \quad (25.30)$$

What a simple equation! Of course, it is necessary to know what the symbol means. U is a physical quantity which we will call the “unworldliness” of the situation. And we have a formula for it. Here is how you calculate the unworldliness. You take all of the known physical laws and write them in a special form. For example, suppose you take the law of mechanics, $F = ma$, and rewrite it as $F - ma = 0$. Then you can call $(F - ma)$ —which should, of course, be zero—the “mismatch,” of mechanics. Next, you take the *square* of this mismatch and call it U_1 , which can be called the “unworldliness of mechanical effects.” In other words, you take

$$U_1 = (F - ma)^2. \quad (25.31)$$

Now you write another physical law, say, $\nabla \cdot E = \rho/\epsilon_0$ and define

$$U_2 = \left(\nabla \cdot E - \frac{\rho}{\epsilon_0} \right)^2,$$

which you might call “the gaussian unworldliness of electricity.” You continue to write U_3, U_4 , and so on—one for every physical law there is

Finally you call the *total* unworldliness U of the world the sum of the various unworldlinesses U_i from all the subphenomena that are involved; that is, $U = \sum U_i$. Then the great “law of nature” is

$$\boxed{U = 0.} \quad (25.32)$$

This “law” means, of course, that the sum of the squares of all the individual mismatches is zero, and the only way the sum of a lot of squares can be zero is for each one of the terms to be zero.

So the “beautifully simple” law in Eq. (25.32) is equivalent to the whole series of equations that you originally wrote down. It is therefore absolutely obvious that a simple notation that just hides the complexity in the definitions of symbols is not real simplicity. *It is just a trick.* The beauty that appears in Eq. (25.32)—just from the fact that several equations are hidden within it—is no more than a trick. When you unwrap the whole thing, you get back where you were before.

However, there *is* more to the simplicity of the laws of electromagnetism written in the form of Eq. (25.29). It means more, just as a theory of vector analysis means more. The fact that the electromagnetic equations can be written in a very particular notation *which was designed* for the four-dimensional geometry of the Lorentz transformations—in other words, as a vector equation in the four-space—means that it is invariant under the Lorentz transformations. It is because the Maxwell equations are invariant under those transformations that they can be written in a beautiful form.

It is no accident that the equations of electrodynamics can be written in the beautifully elegant form of Eq. (25.29). The theory of relativity was developed *because it was found experimentally* that the phenomena predicted by Maxwell’s equations were the same in all inertial systems. And it was precisely by studying the transformation properties of Maxwell’s equations that Lorentz discovered his transformation as the one which left the equations invariant.

There is, however, another reason for writing our equations this way. It has been discovered—after Einstein guessed that it might be so—that *all* of the laws of physics are invariant under the Lorentz transformation. That is the principle of relativity. Therefore, if we invent a notation which shows immediately when a law is written down whether it is invariant or not, we can be sure that in trying to make new theories we will write only equations which are consistent with the principle of relativity.

The fact that the Maxwell equations are simple in this particular notation is not a miracle, because the notation was invented with them in mind. But the interesting physical thing is that *every law* of physics—the propagation of meson waves or the behavior of neutrinos in beta decay, and so forth—must have this same invariance under the same transformation. Then when you are moving at a uniform velocity in a spaceship, all of the laws of nature transform together in such a way that no new phenomenon will show up. It is because the principle of relativity is a fact of nature that in the notation of four-dimensional vectors the equations of the world will look simple.

Lorentz Transformations of the Fields

26-1 The four-potential of a moving charge

We saw in the last chapter that the potential $A_\mu = (\phi, \mathbf{A})$ is a four-vector. The time component is the scalar potential ϕ , and the three space components are the vector potential \mathbf{A} . We also worked out the potentials of a particle moving with uniform speed on a straight line by using the Lorentz transformation (We had already found them by another method in Chapter 21.) For a point charge whose position at the time t is $(vt, 0, 0)$, the potentials at the point (x, y, z) are

$$\begin{aligned}\phi &= \frac{q}{4\pi\epsilon_0\sqrt{1-v^2} \left[\frac{(x-vt)^2}{1-v^2} + y^2 + z^2 \right]^{1/2}} \\ A_x &= \frac{qv}{4\pi\epsilon_0\sqrt{1-v^2} \left[\frac{(x-vt)^2}{1-v^2} + y^2 + z^2 \right]^{1/2}} \\ A_y &= A_z = 0\end{aligned}\quad (26.1)$$

Equations (26.1) give the potentials at x, y , and z at the time t , for a charge whose "present" position (by which we mean the position *at the time t*) is at $x = vt$. Notice that the equations are in terms of $(x - vt)$, y , and z , which are the coordinates measured *from the current position P* of the moving charge (see Fig. 26-1). The actual influence we know really travels at the speed c , so it is the behavior of the charge back at the retarded position P' that really counts.† The point P' is at $x = vt'$ (where, $t' = t - r'/c$ is the retarded time). But we said that the charge was moving with uniform velocity in a straight line, so naturally the behavior at P' and the current position are directly related. In fact, if we make the added assumption that the potentials depend only upon the position and the velocity at the retarded moment, we have in equations (26.1) a *complete* formula for the potentials for a charge moving *any* way. It works this way. Suppose that you have a charge moving in some arbitrary fashion, say with the trajectory in Fig. 26-2, and you are trying to find the potentials at the point (x, y, z) . First, you find the retarded position P' and the velocity v' at that point. Then you imagine that the charge would keep on moving with this velocity during the delay time $(t' - t)$, so that it would then appear at an imaginary position P_{proj} , which we can call the "projected position," and would arrive there with the velocity v' . (Of course, it doesn't do that; its real position at t is at P .) Then the potentials at (x, y, z) are just what equations (26.1) would give for the imaginary charge at the projected position P_{proj} . What we are saying is that since the potentials depend only on what the charge is doing at the *retarded* time, the potentials will be the same whether the charge continued moving at a constant velocity or whether it changed its velocity after t' —that is, after the potentials that were going to appear at (x, y, z) at the time t were already determined.

You know, of course, that the moment that we have the formula for the potentials from a charge moving in any manner whatsoever, we have the complete electrodynamics; we can get the potentials of any charge distribution by super-

† The primes used here to indicate the *retarded* positions and times should not be confused with the primes referring to a Lorentz-transformed frame in the preceding chapter.

26-1 The four-potential of a moving charge

26-2 The fields of a point charge with a constant velocity

26-3 Relativistic transformation of the fields

26-4 The equations of motion in relativistic notation

In this chapter: $c = 1$

Review: Chapter 20, Vol. II, *Solution of Maxwell's Equations in Free Space*

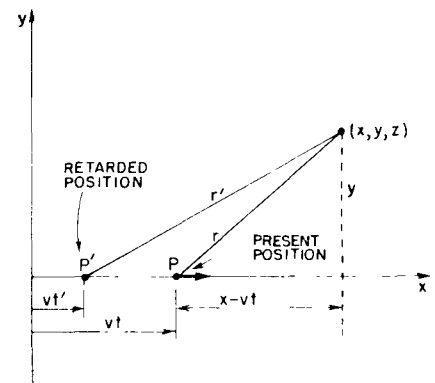


Fig. 26-1. Finding the fields at P due to a charge q moving along the x -axis with the constant speed v . The field "now" at the point (x, y, z) can be expressed in terms of the "present" position P , as well as in terms of P' , the "retarded" position (at $t' = t - r'/c$).

position. Therefore we can summarize all the phenomena of electrodynamics either by writing Maxwell's equations or by the following series of remarks. (Remember them in case you are ever on a desert island. From them, all can be reconstructed. You will, of course, know the Lorentz transformation; you will never forget *that* on a desert island or anywhere else)

First, A_μ is a four-vector. Second, the Coulomb potential for a stationary charge is $q/4\pi\epsilon_0 r$. Third, the potentials produced by a charge moving in any way depend only upon the velocity and position at the retarded time. With those three facts we have everything. From the fact that A_μ is a four-vector, we transform the Coulomb potential, which we know, and get the potentials for a constant velocity. Then, by the last statement that potentials depend only upon the past velocity at the retarded time, we can use the projected position game to find them. It is not a particularly useful way of doing things, but it is interesting to show that the laws of physics can be put in so many different ways.

It is sometimes said, by people who are careless, that all of electrodynamics can be deduced solely from the Lorentz transformation and Coulomb's law. Of course, that is completely false. First, we have to suppose that there is a scalar potential and a vector potential that together make a four-vector. That tells us how the potentials transform. Then why is it that the effects at the retarded time are the only things that count? Better yet, why is it that the potentials depend only on the position and the velocity and not, for instance, on the acceleration? The *fields* \mathbf{E} and \mathbf{B} *do* depend on the acceleration. If you try to make the same kind of an argument with respect to them, you would say that they depend only upon the position and velocity at the retarded time. But then the fields from an accelerating charge would be the same as the fields from a charge at the projected position—which is false. The *fields* depend not only on the position and the velocity along the path but also on the acceleration. So there are several additional tacit assumptions in this great statement that everything can be deduced from the Lorentz transformation. (Whenever you see a sweeping statement that a tremendous amount can come from a very small number of assumptions, you always find that it is false. There are usually a large number of implied assumptions that are far from obvious if you think about them sufficiently carefully.)

26-2 The fields of a point charge with a constant velocity

Now that we have the potentials from a point charge moving at constant velocity, we ought to find the fields—for practical reasons. There are many cases where we have uniformly moving particles—for instance, cosmic rays going through a cloud chamber, or even slow-moving electrons in a wire. So let's at least see what the fields actually do look like for any speed—even for speeds nearly that of light—assuming only that there is no acceleration. It is an interesting question.

We get the fields from the potentials by the usual rules:

$$\mathbf{E} = -\nabla\phi - \frac{\partial \mathbf{A}}{\partial t}, \quad \mathbf{B} = \nabla \times \mathbf{A}.$$

First, for E_z

$$E_z = -\frac{\partial \phi}{\partial z} - \frac{\partial A_z}{\partial t}.$$

But A_z is zero; so differentiating ϕ in equations (26.1), we get

$$E_z = \frac{q}{4\pi\epsilon_0\sqrt{1-v^2}} \left[\frac{z}{1-v^2} + y^2 + z^2 \right]^{3/2}. \quad (26.2)$$

Similarly, for E_y ,

$$E_y = \frac{q}{4\pi\epsilon_0\sqrt{1-v^2}} \left[\frac{(x-vt)^2}{1-v^2} + y^2 + z^2 \right]^{3/2}. \quad (26.3)$$

The x -component is a little more work. The derivative of ϕ is more complicated

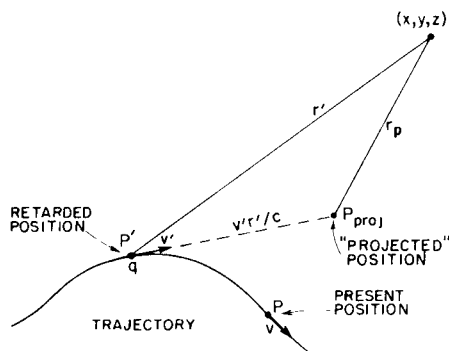


Fig. 26-2. A charge moves on an arbitrary trajectory. The potentials at (x, y, z) at the time t are determined by the position P' and velocity v' at the retarded time $t' = t - r'/c$. They are conveniently expressed in terms of the coordinates from the "projected" position P_{proj} . (The actual position at t is P .)

and A_x is not zero. First,

$$-\frac{\partial \phi}{\partial x} = \frac{q}{4\pi\epsilon_0\sqrt{1-v^2}} \frac{(x-vt)/(1-v^2)}{\left[\frac{(x-vt)^2}{1-v^2} + y^2 + z^2\right]^{3/2}}. \quad (26.4)$$

Then, differentiating A_x with respect to t , we find

$$-\frac{\partial A_x}{\partial t} = \frac{q}{4\pi\epsilon_0\sqrt{1-v^2}} \frac{-v^2(x-vt)/(1-v^2)}{\left[\frac{(x-vt)^2}{1-v^2} + y^2 + z^2\right]^{3/2}}. \quad (26.5)$$

And finally, taking the sum,

$$E_x = \frac{q}{4\pi\epsilon_0\sqrt{1-v^2}} \frac{x-vt}{\left[\frac{(x-vt)^2}{1-v^2} + y^2 + z^2\right]^{3/2}}. \quad (26.6)$$

We'll look at the physics of \mathbf{E} in a minute, let's first find \mathbf{B} . For the z -component,

$$B_z = \frac{\partial A_y}{\partial x} - \frac{\partial A_x}{\partial y}.$$

Since A_y is zero, we have just one derivative to get. Notice, however, that A_x is just $v\phi$, and $\partial/\partial y$ of $v\phi$ is just $-vE_y$. So

$$B_z = vE_y. \quad (26.7)$$

Similarly,

$$B_y = \frac{\partial A_x}{\partial z} - \frac{\partial A_z}{\partial x} = +v \frac{\partial \phi}{\partial z},$$

and

$$B_y = -vE_z. \quad (26.8)$$

Finally, B_x is zero, since A_y and A_z are both zero. We can write the magnetic field simply as

$$\mathbf{B} = \mathbf{v} \times \mathbf{E} \quad (26.9)$$

Now let's see what the fields look like. We will try to draw a picture of the field at various positions around the present position of the charge. It is true that the influence of the charge comes, in a certain sense, from the retarded position, but because the motion is exactly specified, the retarded position is uniquely given in terms of the present position. For uniform velocities, it's nicer to relate the fields to the current position, because the field components at (x, y, z) depend only on $(x - vt)$, y , and z —which are the components of the displacements \mathbf{r}_P from the present position to (x, y, z) (see Fig. 26-3).

Consider first a point with $z = 0$. Then \mathbf{E} has only x - and y -components. From Eqs. (26.3) and (26.6), the ratio of these components is just equal to the ratio of the x - and y -components of the displacement. That means that \mathbf{E} is in the same direction as \mathbf{r}_P , as shown in Fig. 26-3. Since E_z is also proportional to z , it is clear that this result holds in three dimensions. In short, the electric field is radial from the charge, and the field lines radiate directly out of the charge, just as they do for a stationary charge. Of course, the field isn't exactly the same as for the stationary charge, because of all the extra factors of $(1 - v^2)$. But we can show something rather interesting. The difference is just what you would get if you were to draw the Coulomb field with a peculiar set of coordinates in which the scale of x was squashed up by the factor $\sqrt{1 - v^2}$. If you do that, the field lines will be spread out ahead and behind the charge and will be squeezed together around the sides, as shown in Fig. 26-4.

If we relate the strength of \mathbf{E} to the density of the field lines in the conventional way, we see a stronger field at the sides and a weaker field ahead and behind, which is just what the equations say. First, if we look at the strength of the field at right angles to the line of motion, that is, for $(x - vt) = 0$, the distance from

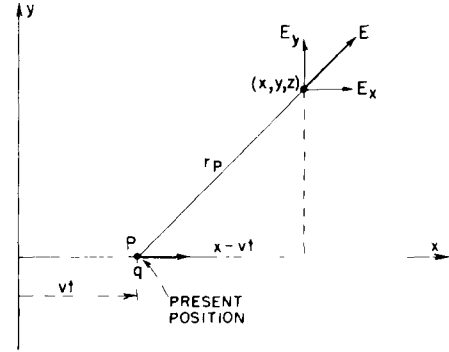


Fig. 26-3. For a charge moving with constant speed, the electric field points radially from the "present" position of the charge.

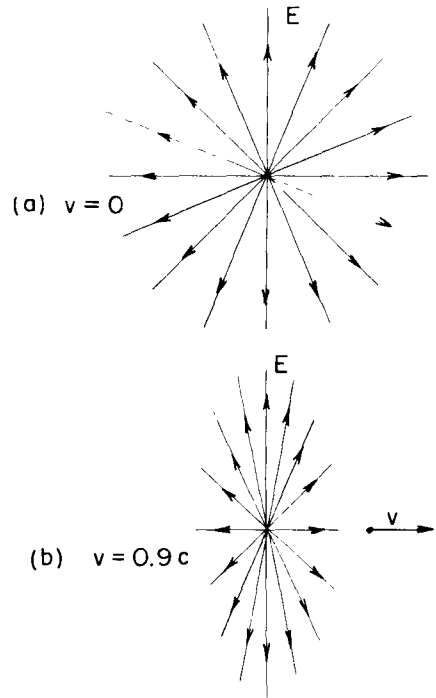


Fig. 26-4. The electric field of a charge moving with the constant speed $v = 0.9c$, part (b), compared with the field of a charge at rest, part (a).

the charge is $(y^2 + z^2)$. Here the total field strength is $\sqrt{E_y^2 + E_z^2}$, which is

$$E = \frac{q}{4\pi\epsilon_0\sqrt{1-v^2}} \frac{1}{y^2 + z^2}. \quad (26.10)$$

The field is proportional to the inverse square of the distance—just like the Coulomb field except increased by the constant, extra factor $1/\sqrt{1-v^2}$, which is always greater than one. So at the *sides* of a moving charge, the electric field is stronger than you get from the Coulomb law. In fact, the field in the sidewise direction is bigger than the Coulomb potential by the ratio of the energy of the particle to its rest mass.

Ahead of the charge (and behind), y and z are zero and

$$E = E_x = \frac{q(1-v^2)}{4\pi\epsilon_0(x-vt)^2}. \quad (26.11)$$

The field again varies as the inverse square of the distance from the charge but is now *reduced* by the factor $(1-v^2)$, in agreement with the picture of the field lines. If v/c is small, v^2/c^2 is still smaller, and the effect of the $(1-v^2)$ terms is very small; we get back to Coulomb's law. But if a particle is moving very close to the speed of light, the field in the forward direction is enormously reduced, and the field in the sidewise direction is enormously increased.

Our results for the electric field of a charge can be put this way: Suppose you were to draw on a piece of paper the field lines for a charge at rest, and then set the picture to travelling with the speed v . Then, of course, the whole picture would be compressed by the Lorentz contraction; that is, the carbon granules on the paper would appear in different places. The miracle of it is that the picture you would see as the page flies by would still represent the field lines of the point charge. The contraction moves them closer together at the sides and spreads them out ahead and behind, just in the right way to give the correct line densities. We have emphasized before that field lines are not real but are only one way of representing the field. However, here they almost seem to be real. In this particular case, if you make the mistake of thinking that the field lines are somehow really there in space, and transform them, you get the correct field. That doesn't, however, make the field lines any more real. All you need do to remind yourself that they aren't real is to think about the electric fields produced by a charge together with a magnet; when the magnet moves, new electric fields are produced, and destroy the beautiful picture. So the neat idea of the contracting picture doesn't work in general. It is, however, a handy way to remember what the fields from a fast-moving charge are like.

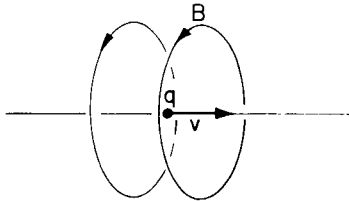


Fig. 26-5. The magnetic field near a moving charge is $\mathbf{v} \times \mathbf{E}$. (Compare with Fig. 26-4.)

The magnetic field is $\mathbf{v} \times \mathbf{E}$ [from Eq. (26.9)]. If you take the velocity crossed into a radial \mathbf{E} -field, you get a \mathbf{B} which circles around the line of motion, as shown in Fig. 26-5. If we put back the c 's, you will see that it's the same result we had for low-velocity charges. A good way to see where the c 's must go is to refer back to the force law,

$$\mathbf{F} = q(\mathbf{E} + \mathbf{v} \times \mathbf{B}).$$

You see that a velocity times the magnetic field has the same dimensions as an electric field. So the right-hand side of Eq. (26.9) must have a factor $1/c^2$:

$$\mathbf{B} = \frac{\mathbf{v} \times \mathbf{E}}{c^2}. \quad (26.12)$$

For a slow-moving charge ($v \ll c$), we can take for \mathbf{E} the Coulomb field; then

$$\mathbf{B} = \frac{q}{4\pi\epsilon_0 c^2} \frac{\mathbf{v} \times \mathbf{r}}{r^3}. \quad (26.13)$$

This formula corresponds exactly to equations for the magnetic field of a current that we found in Section 14-7.

We would like to point out, in passing, something interesting for you to think about. (We will come back to discuss it again later.) Imagine two electrons with velocities at right angles, so that one will cross over the path of the other, but in front of it, so they don't collide. At some instant, their relative positions will be as in Fig 26-6(a). We look at the force on q_1 due to q_2 and vice versa. On q_2 there is only the electric force from q_1 , since q_1 makes no magnetic field along its line of motion. On q_1 , however, there is again the electric force but, in addition, a magnetic force, since it is moving in a \mathbf{B} -field made by q_2 . The forces are as drawn in Fig. 26-6(b). The electric forces on q_1 and q_2 are equal and opposite. However, there is a sidewise (magnetic) force on q_1 and *no sidewise force on q_2* . Does action not equal reaction? We leave it for you to worry about.

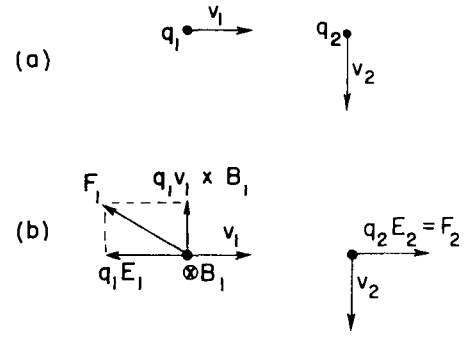


Fig. 26-6. The forces between two moving charges are not always equal and opposite. It appears that "action" is not equal to "reaction."

26-3 Relativistic transformation of the fields

In the last section we calculated the electric and magnetic fields from the transformed potentials. The fields are important, of course, in spite of the arguments given earlier that there is physical meaning and reality to the potentials. The fields, too, are real. It would be convenient for many purposes to have a way to compute the fields in a moving system if you already know the fields in some "rest" system. We have the transformation laws for ϕ and \mathbf{A} , because A_μ is a four-vector. Now we would like to know the transformation laws of \mathbf{E} and \mathbf{B} . Given \mathbf{E} and \mathbf{B} in one frame, how do they look in another frame moving past? It is a convenient transformation to have. We could always work back through the potentials, but it is useful sometimes to be able to transform the fields directly. We will now see how that goes.

How can we find the transformation laws of the fields? We know the transformation laws of the ϕ and \mathbf{A} , and we know how the fields are given in terms of ϕ and \mathbf{A} —it should be easy to find the transformation for the \mathbf{B} and \mathbf{E} . (You might think that with every vector there should be something to make it a four-vector, so with \mathbf{E} there's got to be something else we can use for the fourth component. And also for \mathbf{B} . But it's not so. It's quite different from what you would expect.) To begin with, let's take just a magnetic field \mathbf{B} , which is, of course $\nabla \times \mathbf{A}$. Now we know that the vector potential with its x -, y -, and z -components is only a piece of something; there is also a t -component. Also we know that for derivatives like ∇ , besides the x , y , z parts, there is also a derivative with respect to t . So let's try to figure out what happens if we replace a " y " by a " t ", or a " z " by a " t ," or something like that.

First, notice the form of the terms in $\nabla \times \mathbf{A}$ when we write out the components:

$$B_x = \frac{\partial A_z}{\partial y} - \frac{\partial A_y}{\partial z}, \quad B_y = \frac{\partial A_x}{\partial z} - \frac{\partial A_z}{\partial x}, \quad B_z = \frac{\partial A_y}{\partial x} - \frac{\partial A_x}{\partial y}. \quad (26.14)$$

The x -component is equal to a couple of terms that involve only y - and z -components. Suppose we call this combination of derivatives and components a "zy-thing," and give it a shorthand name, F_{zy} . We simply mean that

$$F_{zy} \equiv \frac{\partial A_z}{\partial y} - \frac{\partial A_y}{\partial z}. \quad (26.15)$$

Similarly, B_y is equal to the same kind of "thing," but this time it is an "xz-thing." And B_z is, of course, the corresponding "yx-thing." We have

$$B_x = F_{zy}, \quad B_y = F_{xz}, \quad B_z = F_{yx}. \quad (26.16)$$

Now what happens if we simply try to concoct also some " t "-type things like F_{xt} and F_{tz} (since nature should be nice and symmetric in x , y , z , and t)? For instance, what is F_{tz} ? It is, of course,

$$\frac{\partial A_t}{\partial z} - \frac{\partial A_z}{\partial t}.$$

But remember that $A_t = \phi$, so it is also

$$\frac{\partial \phi}{\partial z} = \frac{\partial A_z}{\partial t}.$$

You've seen that before. It is the z -component of \mathbf{E} . Well, almost—there is a sign wrong. But we forgot that in the four-dimensional gradient the t -derivative comes with the opposite sign from x , y , and z . So we should really have taken the more consistent extension of F_{tz} as

$$F_{tz} = \frac{\partial A_t}{\partial z} + \frac{\partial A_z}{\partial t}. \quad (26.17)$$

Then it is exactly equal to $-E_z$. Trying also F_{tx} and F_{ty} , we find that the three possibilities give

$$F_{tx} = -E_x, \quad F_{ty} = -E_y, \quad F_{tz} = -E_z. \quad (26.18)$$

What happens if both subscripts are t ? Or, for that matter, if both are x ? We get things like

$$F_{tt} = \frac{\partial A_t}{\partial t} - \frac{\partial A_t}{\partial t},$$

and

$$F_{xx} = \frac{\partial A_x}{\partial x} - \frac{\partial A_x}{\partial x},$$

which give nothing but zero.

We have then six of these F -things. There are six more which you get by reversing the subscripts, but they give nothing really new, since

$$F_{yx} = -F_{xy}.$$

and so on. So, out of sixteen possible combinations of the four subscripts taken in pairs, we get only six different physical objects; *and they are the components of \mathbf{B} and \mathbf{E} .*

To represent the general term of F , we will use the general subscripts μ and ν , where each can stand for 0, 1, 2, or 3—meaning in our usual four-vector notation t , x , y , and z . Also, everything will be consistent with our four-vector notation if we define $F_{\mu\nu}$ by

$$F_{\mu\nu} = \nabla_\mu A_\nu - \nabla_\nu A_\mu, \quad (26.19)$$

remembering that $\nabla_\mu = (\partial/\partial t, -\partial/\partial x, -\partial/\partial y, -\partial/\partial z)$ and that $A_\mu = (\phi, A_x, A_y, A_z)$.

What we have found is that there are six quantities that belong together in nature—that are different aspects of the same thing. The electric and magnetic fields which we have considered as separate vectors in our slow-moving world (where we don't worry about the speed of light) are not vectors in four-space. They are parts of a new "thing." Our physical "field" is really the six-component object $F_{\mu\nu}$. That is the way we must look at it for relativity. We summarize our results on $F_{\mu\nu}$ in Table 26-1.

You see that what we have done here is to generalize the cross product. We began with the curl operation, and the fact that the transformation properties of the curl are the same as the transformation properties of *two* vectors—the ordinary three-dimensional vector \mathbf{A} and the gradient operator which we know also behaves like a vector. Let's look for a moment at an ordinary cross product in three dimensions, for example, the angular momentum of a particle. When an object is moving in a plane, the quantity $(xv_y - yv_x)$ is important. For motion in three dimensions, there are three such important quantities, which we call the angular momentum:

$$L_{xy} = m(xv_y - yv_x), \quad L_{yz} = m(yv_z - zv_y), \quad L_{zx} = m(zv_x - xv_z).$$

Then (although you may have forgotten by now) we discovered in Chapter 20 of Vol. I the miracle that these three quantities could be identified with the com-

Table 26-1

The components of $F_{\mu\nu}$

$F_{\mu\nu} = -F_{\nu\mu}$	
$F_{\mu\mu} = 0$	
$F_{xy} = -B_z$	$F_{xt} = E_x$
$F_{yz} = -B_x$	$F_{yt} = E_y$
$F_{zx} = -B_y$	$F_{zt} = E_z$

ponents of a vector. In order to do so, we had to make an artificial rule with a right-hand convention. It was just luck. It was luck because L_{ij} (with i and j equal to x, y , or z) was an antisymmetric object

$$L_{ij} = -L_{ji}, \quad L_{ii} = 0$$

Of the nine possible quantities, there are only three independent numbers. And it just happens that when you change coordinate systems these three objects transform in exactly the same way as the components of a vector.

The same thing lets us represent an element of surface as a vector. A surface element has two parts—say dx and dy —which we can represent by the vector $d\mathbf{a}$ normal to the surface. But we can't do that in four dimensions. What is the “normal” to $dx\,dy$? Is it along z or along t ?

In short, for three dimensions it happens by luck that after you've taken a combination of two vectors like L_{ij} , you can represent it again by another vector because there are just three terms that happen to transform like the components of a vector. But in four dimensions that is evidently impossible, because there are six independent terms, and you can't represent six things by four things.

Even in three dimensions it is possible to have combinations of vectors that can't be represented by vectors. Suppose we take any two vectors $\mathbf{a} = (a_x, a_y, a_z)$ and $\mathbf{b} = (b_x, b_y, b_z)$, and make the various possible combinations of components, like $a_x b_x, a_x b_y$, etc. There would be nine possible quantities:

$$\begin{array}{lll} a_x b_x, & a_x b_y, & a_x b_z, \\ a_y b_x, & a_y b_y, & a_y b_z, \\ a_z b_x, & a_z b_y, & a_z b_z. \end{array}$$

We might call these quantities T_{ij} .

If we now go to a rotated coordinate system (say rotated about the z -axis), the components of \mathbf{a} and \mathbf{b} are changed. In the new system, a_x , for example, gets replaced by

$$a'_x = a_x \cos \theta + a_y \sin \theta,$$

and b_y gets replaced by

$$b'_y = b_y \cos \theta - b_x \sin \theta.$$

And similarly for other components. The nine components of the product quantity T_{ij} we have invented are all changed too, of course. For instance, $T_{xy} = a_x b_y$ gets changed to

$$T'_{xy} = a_x b_y (\cos^2 \theta) - a_x b_x (\cos \theta \sin \theta) + a_y b_y (\sin \theta \cos \theta) - a_y b_x (\sin^2 \theta),$$

or

$$T'_{xy} = T_{xy} \cos^2 \theta - T_{xx} \cos \theta \sin \theta + T_{yy} \sin \theta \cos \theta - T_{yx} \sin^2 \theta.$$

Each component of T'_{ij} is a linear combination of the components of T_{ij} .

So we discover that it is not only possible to have a “vector product” like $\mathbf{a} \times \mathbf{b}$ which has three components that transform like a vector, but we can—artificially—also make another kind of “product” of two vectors T_{ij} with nine components that transform under a rotation by a complicated set of rules that we could figure out. Such an object which has two indices to describe it, instead of one, is called a *tensor*. It is a tensor of the “second rank,” because you can play this game with three vectors too and get a tensor of the third rank,—or with four, to get a tensor of the fourth rank, and so on. A tensor of the first rank is a vector.

The point of all this is that our electromagnetic quantity $F_{\mu\nu}$ is also a tensor of the second rank, because it has two indices in it. It is, however, a tensor in four dimensions. It transforms in a special way which we will work out in a moment—it is just the way a product of vectors transforms. For $F_{\mu\nu}$ it happens that if you change the indices around, $F_{\mu\nu}$ changes sign. That's a special case—it is

an *antisymmetric tensor*. So we say: the electric and magnetic fields are both part of an antisymmetric tensor of the second rank in four dimensions.

You've come a long way. Remember way back when we defined what a velocity meant? Now we are talking about "an antisymmetric tensor of the second rank in four dimensions."

Now we have to find the law of the transformation of $F_{\mu\nu}$. It isn't at all difficult to do; it's just laborious—the brains involved are nil, but the work is not. What we want is the Lorentz transformation of $\nabla_\mu A_\nu - \nabla_\nu A_\mu$. Since ∇_μ is just a special case of a vector, we will work with the general antisymmetric vector combination, which we can call $G_{\mu\nu}$:

$$G_{\mu\nu} = a_\mu b_\nu - a_\nu b_\mu. \quad (26.20)$$

(For our purposes, a_μ will eventually be replaced by ∇_μ and b_μ will be replaced by the potential A_μ .) The components of a_μ and b_μ transform by the Lorentz formulas, which are

$$\begin{aligned} a'_t &= \frac{a_t - va_x}{\sqrt{1-v^2}}, & b'_t &= \frac{b_t - vb_x}{\sqrt{1-v^2}}, \\ a'_x &= \frac{a_x - va_t}{\sqrt{1-v^2}}, & b'_x &= \frac{b_x - vb_t}{\sqrt{1-v^2}}, \\ a'_y &= a_y, & b'_y &= b_y, \\ a'_z &= a_z, & b'_z &= b_z. \end{aligned} \quad (26.21)$$

Now let's transform the components of $G_{\mu\nu}$. We start with G_{tx} :

$$\begin{aligned} G'_{tx} &= a'_t b'_x - a'_x b'_t \\ &= \left(\frac{a_t - va_x}{\sqrt{1-v^2}} \right) \left(\frac{b_x - vb_t}{\sqrt{1-v^2}} \right) - \left(\frac{a_x - va_t}{\sqrt{1-v^2}} \right) \left(\frac{b_t - vb_x}{\sqrt{1-v^2}} \right) \\ &= a_t b_x - a_x b_t. \end{aligned}$$

But that is just G_{tx} ; so we have the simple result

$$G'_{tx} = G_{tx}.$$

We will do one more

$$G'_{ty} = \frac{a_t - va_x}{\sqrt{1-v^2}} b_y - a_y \frac{b_t - vb_x}{\sqrt{1-v^2}} = \frac{(a_t b_y - a_y b_t) - v(a_x b_y - a_y b_x)}{\sqrt{1-v^2}}.$$

So we get that

$$G'_{ty} = \frac{G_{ty} - vG_{xy}}{\sqrt{1-v^2}}.$$

And, of course, in the same way,

$$G'_{tz} = \frac{G_{tz} - vG_{xz}}{\sqrt{1-v^2}}.$$

It is clear how the rest will go. Let's make a table of all six terms, only now we may as well write them for $F_{\mu\nu}$:

$$\begin{aligned} F'_{tx} &= F_{tx}, & F'_{xy} &= \frac{F_{xy} - vF_{ty}}{\sqrt{1-v^2}}, \\ F'_{ty} &= \frac{F_{ty} - vF_{xy}}{\sqrt{1-v^2}}, & F'_{yz} &= F_{yz}, \\ F'_{tz} &= \frac{F_{tz} - vF_{xz}}{\sqrt{1-v^2}}, & F'_{zx} &= \frac{F_{zx} - vF_{zt}}{\sqrt{1-v^2}}. \end{aligned} \quad (26.22)$$

Of course, we still have $F'_{\mu\nu} = -F'_{\nu\mu}$ and $F_{\mu\mu} = 0$.

So we have the transformation of the electric and magnetic fields. All we have to do is look at Table 26-1 to find out what our grand notation in terms of $F_{\mu\nu}$ means in terms of \mathbf{E} and \mathbf{B} . It's just a matter of substitution. So that we can see how it looks in the ordinary symbols, we'll rewrite our transformation of the field components in Table 26-2.

Table 26-2

The Lorentz transformation of the electric and magnetic fields (Note: $c = 1$)

$E'_x = E_x$	$B'_x = B_x$
$E'_y = \frac{E_y - vB_z}{\sqrt{1 - v^2}}$	$B'_y = \frac{B_y + vE_z}{\sqrt{1 - v^2}}$
$E'_z = \frac{E_z + vB_y}{\sqrt{1 - v^2}}$	$B'_z = \frac{B_z - vE_y}{\sqrt{1 - v^2}}$

The equations in Table 26-2 tell us how \mathbf{E} and \mathbf{B} change if we go from one inertial frame to another. If we know \mathbf{E} and \mathbf{B} in one system, we can find what they are in another that moves by with the speed v .

We can write these equations in a form that is easier to remember if we notice that since v is in the x -direction, all the terms with v are components of the cross products $\mathbf{v} \times \mathbf{E}$ and $\mathbf{v} \times \mathbf{B}$. So we can rewrite the transformations as shown in Table 26-3

Table 26-3

An alternative form for the field transformations (Note: $c = 1$)

$E'_x = E_x$	$B'_x = B_x$
$E'_y = \frac{(E + \mathbf{v} \times \mathbf{B})_y}{\sqrt{1 - v^2}}$	$B'_y = \frac{(\mathbf{B} - \mathbf{v} \times \mathbf{E})_y}{\sqrt{1 - v^2}}$
$E'_z = \frac{(E + \mathbf{v} \times \mathbf{B})_z}{\sqrt{1 - v^2}}$	$B'_z = \frac{(\mathbf{B} - \mathbf{v} \times \mathbf{E})_z}{\sqrt{1 - v^2}}$

It is now easier to remember which components go where. In fact, the transformation can be written even more simply if we define the field components along x as the “parallel” components E_{\parallel} and B_{\parallel} (because they are parallel to the relative velocity of S and S'), and the total transverse components—the vector sums of the y - and z -components—as the “perpendicular” components E_{\perp} and B_{\perp} . Then we get the equations in Table 26-4. (We have also put back the c 's, so it will be more convenient when we want to refer back later.)

Table 26-4

Still another form for the Lorentz transformation of \mathbf{E} and \mathbf{B}

$E'_{\parallel} = E_{\parallel}$	$B'_{\parallel} = B_{\parallel}$
$E'_{\perp} = \frac{(\mathbf{E} + \mathbf{v} \times \mathbf{B})_{\perp}}{\sqrt{1 - v^2/c^2}}$	$B'_{\perp} = \frac{(\mathbf{B} - \frac{\mathbf{v} \times \mathbf{E}}{c^2})_{\perp}}{\sqrt{1 - v^2/c^2}}$

The field transformations give us another way of solving some problems we have done before—for instance, for finding the fields of a moving point charge. We have worked out the fields before by differentiating the potentials. But we could now do it by transforming the Coulomb field. If we have a point charge at rest in the S -frame, then there is only the simple radial \mathbf{E} -field. In the S' -frame we will see a point charge moving with the velocity u , if the S' -frame moves by the

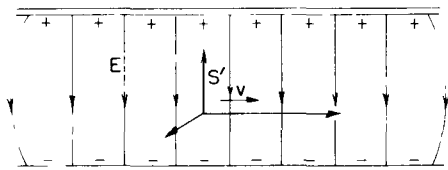


Fig. 26-7. The coordinate frame S' moving through a static electric field.

S -frame with the speed $v = -u$. We will let you show that the transformations of Tables 26-3 and 26-4 give the same electric and magnetic fields we got in Section 26-2.

The transformation of Table 26-2 gives us an interesting and simple answer for what we see if we move past *any* system of fixed charges. For example, suppose we want to know the fields in *our* frame S' if we are moving along between the plates of a condenser, as shown in Fig. 26-7. (It is, of course, the same thing if we say that a charged condenser is moving past *us*.) What do we see? The transformation is easy in this case because the \mathbf{B} -field in the original system is zero. Suppose, first, that our motion is perpendicular to \mathbf{E} , then we will see an $\mathbf{E}' = \mathbf{E}/\sqrt{1 - v^2/c^2}$ which is still completely transverse. We will see, in addition, a magnetic field $\mathbf{B}' = -\mathbf{v} \times \mathbf{E}'/c^2$. (The $\sqrt{1 - v^2/c^2}$ doesn't appear in our formula for \mathbf{B}' because we wrote it in terms of \mathbf{E}' rather than \mathbf{E} ; but it's the same thing.) So when we move along perpendicular to a static electric field, we see a reduced \mathbf{E} and an added transverse \mathbf{B} . If our motion is not perpendicular to \mathbf{E} , we break \mathbf{E} into \mathbf{E}_{\parallel} and \mathbf{E}_{\perp} . The parallel part is unchanged, $E'_{\parallel} = E_{\parallel}$, and the perpendicular component does as just described.

Let's take the opposite case, and imagine we are moving through a pure static *magnetic* field. This time we would see an *electric* field \mathbf{E}' equal to $\mathbf{v} \times \mathbf{B}'$, and the magnetic field changed by the factor $1/\sqrt{1 - v^2/c^2}$ (assuming it is transverse). So long as v is much less than c , we can neglect the change in the magnetic field, and the main effect is that an electric field appears. As one example of this effect, consider this once famous problem of determining the speed of an airplane. It's no longer famous, since radar can now be used to determine the air speed from ground reflections, but for many years it was very hard to find the speed of an airplane in bad weather. You could not see the ground and you didn't know which way was up, and so on. Yet it was important to know how fast you were moving relative to the earth. How can this be done without seeing the earth? Many who knew the transformation formulas thought of the idea of using the fact that the airplane moves in the magnetic field of the earth. Suppose that an airplane is flying where there is a magnetic field more or less known. Let's just take the simple case where the magnetic field is vertical. If we were flying through it with a horizontal velocity v , then, according to our formula, we should see an electric field which is $\mathbf{v} \times \mathbf{B}$, i.e., perpendicular to the line of motion. If we hang an insulated wire across the airplane, this electric field will induce charges on the ends of the wire. That is nothing new. From the point of view of someone on the ground, we are moving a wire through a field, and the $\mathbf{v} \times \mathbf{B}$ force causes charges to move to the ends of the wire. The transformation equations just say the same thing in a different way. (The fact that we can say the thing more than one way doesn't mean that one way is better than another. We are getting so many different methods and tools that we can usually get the same result in 65 different ways!)

So to measure v , all we have to do is measure the voltage between the ends of the wire. We can't do it with a voltmeter because the same fields will act on the wires in the voltmeter, but there are ways of measuring such fields. We talked about some of them when we discussed atmospheric electricity in Chapter 9. So it should be possible to measure the speed of the airplane.

This important problem was, however, never solved this way. The reason is that the electric field that is developed is of the order of millivolts per meter. It is possible to measure such fields, but the trouble is that these fields are, unfortunately, not any different from any other electric fields. The field that is produced by motion through the magnetic field can't be distinguished from some electric field that was already in the air from another cause, say from electrostatic charges in the air, or on the clouds. We described in Chapter 9 that there are, typically, electric fields above the surface of the earth with strengths of about 100 volts per meter. But they are quite irregular. So as the airplane flies through the air, it sees fluctuations of atmospheric electric fields which are enormous in comparison to the tiny fields produced by the $\mathbf{v} \times \mathbf{B}$ term, and it turns out for practical reasons to be impossible to measure speeds of an airplane by its motion through the earth's magnetic field.

26-4 The equations of motion in relativistic notation*

It doesn't do much good to find electric and magnetic fields from Maxwell's equations unless we know what the fields do when we have them. You may remember that the fields are required to find the forces on charges, and that those forces determine the motion of the charge. So, of course, part of the theory of electrodynamics is the relation between the motion of charges and the forces.

For a single charge in the fields \mathbf{E} and \mathbf{B} , the force is

$$\mathbf{F} = q(\mathbf{E} + \mathbf{v} \times \mathbf{B}). \quad (26.23)$$

This force is equal to the mass times the acceleration for low velocities, but the correct law for any velocity is that the force is equal to $d\mathbf{p}/dt$. Writing $\mathbf{p} = m_0\mathbf{v}/\sqrt{1 - v^2/c^2}$, we find that the relativistically correct equation of motion is

$$\frac{d}{dt} \left(\frac{m_0\mathbf{v}}{\sqrt{1 - v^2/c^2}} \right) = \mathbf{F} = q(\mathbf{E} + \mathbf{v} \times \mathbf{B}). \quad (26.24)$$

We would like now to discuss this equation from the point of view of relativity. Since we have put our Maxwell equations in relativistic form, it would be interesting to see what the equations of motion would look like in relativistic form. Let's see whether we can rewrite the equation in a four-vector notation.

We know that the momentum is part of a four-vector p_μ whose time component is the energy $m_0c/\sqrt{1 - v^2/c^2}$. So we might think to replace the left-hand side of Eq. (26.24) by dp_μ/dt . Then we need only find a fourth component to go with \mathbf{F} . This fourth component must equal the rate-of-change of the energy, or the rate of doing work, which is $\mathbf{F} \cdot \mathbf{v}$. We would then like to write the right-hand side of Eq. (26.24) as a four-vector like $(\mathbf{F} \cdot \mathbf{v}, F_x, F_y, F_z)$. But this does not make a four-vector.

The *time* derivative of a four-vector is no longer a four-vector, because the d/dt requires the choice of some special frame for measuring t . We got into that trouble before when we tried to make \mathbf{v} into a four-vector. Our first guess was that the time component would be $cdt/dt = c$. But the quantities

$$\left(c, \frac{dx}{dt}, \frac{dy}{dt}, \frac{dz}{dt} \right) = (c, \mathbf{v}) \quad (26.25)$$

are *not* the components of a four-vector. We found that they could be made into one by multiplying each component by $1/\sqrt{1 - v^2/c^2}$. The "four-velocity" u_μ is the four-vector

$$u_\mu = \left(\frac{c}{\sqrt{1 - v^2/c^2}}, \frac{\mathbf{v}}{\sqrt{1 - v^2/c^2}} \right). \quad (26.26)$$

So it appears that the trick is to multiply d/dt by $1/\sqrt{1 - v^2/c^2}$, if we want the derivatives to make a four-vector.

Our second guess then is that

$$\frac{1}{\sqrt{1 - v^2/c^2}} \frac{d}{dt} (p_\mu) \quad (26.27)$$

should be a four-vector. But what is \mathbf{v} ? It is the velocity of the particle—not of a coordinate frame! Then the quantity f_μ defined by

$$f_\mu = \left(\frac{\mathbf{F} \cdot \mathbf{v}}{\sqrt{1 - v^2/c^2}}, \frac{\mathbf{F}}{\sqrt{1 - v^2/c^2}} \right) \quad (26.28)$$

is the extension into four dimensions of a force—we can call it the "four-force." It is indeed a four-vector, and its space components are not the components of \mathbf{F} but of $\mathbf{F}/\sqrt{1 - v^2/c^2}$.

* In this section we will put back all of the c 's.

The question is—why is f_μ a four-vector? It would be nice to get a little understanding of that $1/\sqrt{1 - v^2/c^2}$ factor. Since it has come up twice now, it is time to see why the d/dt can always be fixed by the same factor. The answer is in the following: When we take the time derivative of some function x , we compute the increment Δx in a small interval Δt in the variable t . But in another frame, the interval Δt might correspond to a change in both t' and x' , so if we vary only t' , the change in x will be different. We have to find a variable for our differentiation that is a measure of an “interval” in *space-time*, which will then be the same in all coordinate systems. When we take Δx for that interval, it will be the same for all coordinate frames. When a particle “moves” in four-space, there are the changes $\Delta t, \Delta x, \Delta y, \Delta z$. Can we make an invariant interval out of them? Well, they are the components of the four-vector $x_\mu = (ct, x, y, z)$ so if we define a quantity Δs by

$$(\Delta s)^2 = \frac{1}{c^2} \Delta x_\mu \Delta x_\mu = \frac{1}{c^2} (c^2 \Delta t^2 - \Delta x^2 - \Delta y^2 - \Delta z^2) \quad (26.29)$$

—which is a four-dimensional dot product—we then have a good four-scalar to use as a measure of a four-dimensional interval. From Δs —or its limit ds —we can define a parameter $s = \int ds$. And a derivative with respect to s , d/ds , is a nice four-dimensional operation, because it is invariant with respect to a Lorentz transformation.

It is easy to relate ds to dt for a moving particle. For a moving point particle,

$$dx = v_x dt, \quad dy = v_y dt, \quad dz = v_z dt, \quad (26.30)$$

and

$$ds = \sqrt{(dt^2/c^2)(c^2 - v_x^2 - v_y^2 - v_z^2)} = dt\sqrt{1 - v^2/c^2}. \quad (26.31)$$

So the operator

$$\frac{1}{\sqrt{1 - v^2/c^2}} \frac{d}{dt}$$

is an *invariant operator*. If we operate on any four-vector with it, we get another four-vector. For instance, if we operate on (ct, x, y, z) , we get the four-velocity u_μ :

$$\frac{dx_\mu}{ds} = u_\mu.$$

We see now why the factor $\sqrt{1 - v^2/c^2}$ fixes things up.

The invariant variable s is a useful physical quantity. It is called the “proper time” along the path of a particle, because ds is always an interval of time in a frame that is moving with the particle at any particular instant. (Then, $\Delta x = \Delta y = \Delta z = 0$, and $\Delta s = \Delta t$.) If you can imagine some “clock” whose rate doesn’t depend on the acceleration, such a clock carried along with the particle would show the time s .

We can now go back and write Newton’s law (as corrected by Einstein) in the neat form

$$\frac{dp_\mu}{ds} = f_\mu, \quad (26.32)$$

where f_μ is given in Eq. (26.28). Also, the momentum p_μ can be written as

$$p_\mu = m_0 u_\mu = m_0 \frac{dx_\mu}{ds}, \quad (26.33)$$

where the coordinates $x_\mu = (ct, x, y, z)$ now describe the trajectory of the particle. Finally, the four-dimensional notation gives us this very simple form of the equations of motion:

$$f_\mu = m_0 \frac{d^2 x_\mu}{ds^2}, \quad (26.34)$$

which is reminiscent of $F = ma$. It is important to notice that Eq. (26.34) is *not* the same as $F = ma$, because the four-vector formula Eq. (26.34) has in it the

relativistic mechanics which are different from Newton's law for high velocities. It is unlike the case of Maxwell's equations, where we were able to rewrite the equations in the relativistic form *without any change in the meaning at all*—but with just a change of notation.

Now let's return to Eq. (26.24) and see how we can write the right-hand side in four-vector notation. The three components—when divided by $\sqrt{1 - v^2/c^2}$ —are the components of f_μ , so

$$f_x = \frac{q(\mathbf{E} + \mathbf{v} \times \mathbf{B})_x}{\sqrt{1 - v^2/c^2}} = q \left[\frac{E_x}{\sqrt{1 - v^2/c^2}} + \frac{v_y B_z}{\sqrt{1 - v^2/c^2}} - \frac{v_z B_y}{\sqrt{1 - v^2/c^2}} \right]. \quad (26.35)$$

Now we must put all quantities in their relativistic notation. First, $c/\sqrt{1 - v^2/c^2}$ and $v_y/\sqrt{1 - v^2/c^2}$ and $v_z/\sqrt{1 - v^2/c^2}$ are the t -, y -, and z -components of the four-velocity u_μ . And the components of \mathbf{E} and \mathbf{B} are components of the second-rank tensor of the fields $F_{\mu\nu}$. Looking back in Table 26-1 for the components of $F_{\mu\nu}$ that correspond to E_x , B_z , and B_y , we get

$$f_x = q(u_t F_{xt} - u_y F_{xy} - u_z F_{xz}),$$

which begins to look interesting. Every term has the subscript x , which is reasonable, since we're finding an x -component. Then all the others appear in pairs: tt , yy , zz —except that the xx -term is missing. So we just stick it in, and write

$$f_x = q(u_t F_{xt} - u_x F_{xx} - u_y F_{xy} - u_z F_{xz}). \quad (26.36)$$

We haven't changed anything because $F_{\mu\nu}$ is antisymmetric, and F_{xx} is zero. The reason for wanting to put in the xx -term is so that we can write Eq. (26.36) in the short-hand form

$$f_\mu = q u_\nu F_{\mu\nu}. \quad (26.37)$$

This equation is the same as Eq. (26.36) if we make the *rule* that whenever any subscript occurs *twice* (as ν does here), you automatically sum over terms in the same way as for the scalar product, *using the same convention for the signs*.

You can easily believe that (26.37) works equally well for $\mu = y$ or $\mu = z$, but what about $\mu = t$? Let's see, for fun, what it says:

$$f_t = q(u_t F_{tt} - u_x F_{tx} - u_y F_{ty} - u_z F_{tz})$$

Now we have to translate back to E 's and B 's. We get

$$f_t = q \left(0 + \frac{v_x}{\sqrt{1 - v^2/c^2}} E_x + \frac{v_y}{\sqrt{1 - v^2/c^2}} E_y + \frac{v_z}{\sqrt{1 - v^2/c^2}} E_z \right), \quad (26.38)$$

or

$$f_t = \frac{q \mathbf{v} \cdot \mathbf{E}}{\sqrt{1 - v^2/c^2}}$$

But from Eq. (26.28), f_t is supposed to be

$$\frac{\mathbf{F} \cdot \mathbf{v}}{\sqrt{1 - v^2/c^2}} = \frac{q(\mathbf{E} + \mathbf{v} \times \mathbf{B}) \cdot \mathbf{v}}{\sqrt{1 - v^2/c^2}}.$$

This is the same thing as Eq. (26.38), since $(\mathbf{v} \times \mathbf{B}) \cdot \mathbf{v}$ is zero. So everything comes out all right.

Summarizing, our equation of motion can be written in the elegant form

$$m_0 \frac{d^2 x_\mu}{ds^2} = f_\mu = q u_\nu F_{\mu\nu}. \quad (26.39)$$

Although it is nice to see that the equations can be written that way, this form is not particularly useful. It's usually more convenient to solve for particle motions by using the original equations (26.24), and that's what we will usually do.

Field Energy and Field Momentum

27-1 Local conservation

It is clear that the energy of matter is not conserved. When an object radiates light it loses energy. However, the energy lost is possibly describable in some other form, say in the light. Therefore the theory of the conservation of energy is incomplete without a consideration of the energy which is associated with the light or, in general, with the electromagnetic field. We take up now the law of conservation of energy and, also, of momentum for the fields. Certainly, we cannot treat one without the other, because in the relativity theory they are different aspects of the same four-vector.

Very early in Volume I, we discussed the conservation of energy; we said then merely that the total energy in the world is constant. Now we want to extend the idea of the energy conservation law in an important way—in a way that says something in *detail* about *how* energy is conserved. The new law will say that if energy goes away from a region, it is because it *flows* away through the boundaries of that region. It is a somewhat stronger law than the conservation of energy without such a restriction.

To see what the statement means, let's look at how the law of the conservation of charge works. We described the conservation of charge by saying that there is a current density \mathbf{j} and a charge density ρ , and that when the charge decreases at some place there must be a flow of charge away from that place. We call that the conservation of charge. The mathematical form of the conservation law is

$$\nabla \cdot \mathbf{j} = -\frac{\partial \rho}{\partial t}. \quad (27.1)$$

This law has the consequence that the total charge in the world is always constant—there is never any net gain or loss of charge. However, the total charge in the world could be constant in another way. Suppose that there is some charge Q_1 near some point (1) while there is no charge near some point (2) some distance away (Fig. 27-1). Now suppose that, as time goes on, the charge Q_1 were to gradually fade away and that *simultaneously* with the decrease of Q_1 some charge Q_2 would appear near point (2), and in such a way that at every instant the sum of Q_1 and Q_2 was a constant. In other words, at any intermediate state the amount of charge lost by Q_1 would be added to Q_2 . Then the total amount of charge in the world would be conserved. That's a “world-wide” conservation, but not what we will call a “local” conservation, because in order for the charge to get from (1) to (2), it didn't have to appear anywhere in the space between point (1) and point (2). Locally, the charge was just “lost.”

There is a difficulty with such a “world-wide” conservation law in the theory of relativity. The concept of “simultaneous moments” at distant points is one which is not equivalent in different systems. Two events that are simultaneous in one system are not simultaneous for another system moving past. For “world-wide” conservation of the kind described, it is necessary that the charge lost from Q_1 should appear *simultaneously* in Q_2 . Otherwise there would be some moments when the charge was not conserved. There seems to be no way to make the law of charge conservation relativistically invariant without making it a “local” conservation law. As a matter of fact, the requirement of the Lorentz relativistic invariance seems to restrict the possible laws of nature in surprising ways. In modern quantum field theory, for example, people have often wanted to alter the theory by allowing what we call a “nonlocal” interaction—where something *here*

27-1 Local conservation

27-2 Energy conservation and electromagnetism

27-3 Energy density and energy flow in the electromagnetic field

27-4 The ambiguity of the field energy

27-5 Examples of energy flow

27-6 Field momentum

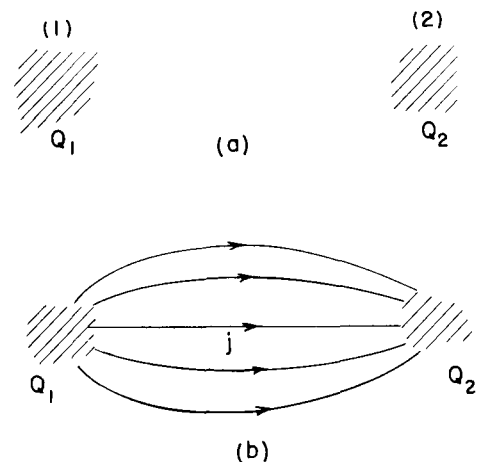


Fig. 27-1. Two ways to conserve charge: (a) $Q_1 + Q_2$ is constant; (b) $dQ_1/dt = \int \mathbf{j} \cdot \mathbf{n} da = -dQ_2/dt$.

has a direct effect on something *there*—but we get in trouble with the relativity principle.

“Local” conservation involves another idea. It says that a charge can get from one place to another only if there is something happening in the space between. To describe the law we need not only the density of charge, ρ , but also another kind of quantity, namely \mathbf{j} , a vector giving the rate of flow of charge across a surface. Then the flow is related to the rate of change of the density by Eq. (27.1). This is the more extreme kind of a conservation law. It says that charge is conserved in a special way—conserved “locally.”

It turns out that energy conservation is also a *local* process. There is not only an energy density in a given region of space but also a vector to represent the rate of flow of the energy through a surface. For example, when a light source radiates, we can find the light energy moving out from the source. If we imagine some mathematical surface surrounding the light source, the energy lost from inside the surface is equal to the energy that flows out through the surface.

27-2 Energy conservation and electromagnetism

We want now to write quantitatively the conservation of energy for electromagnetism. To do that, we have to describe how much energy there is in any volume element of space, and also the rate of energy flow. Suppose we think first only of the electromagnetic field energy. We will let u represent the *energy density* in the field (that is, the amount of energy per unit volume in space) and let the vector \mathbf{S} represent the *energy flux* of the field (that is, the flow of energy per unit time across a unit area perpendicular to the flow). Then, in perfect analogy with the conservation of charge, Eq. (27.1), we can write the “local” law of energy conservation in the field as

$$\frac{\partial u}{\partial t} = -\nabla \cdot \mathbf{S}. \quad (27.2)$$

Of course, this law is not true in general; it is not true that the field energy is conserved. Suppose you are in a dark room and then turn on the light switch. All of a sudden the room is full of light, so there is energy in the field, although there wasn’t any energy there before. Equation (27.2) is not the complete conservation law, because the *field energy alone* is not conserved, only the total energy in the world—there is also the energy of matter. The field energy will change if there is some work being done by matter on the field or by the field on matter.

However, if there is matter inside the volume of interest, we know how much energy it has: Each particle has the energy $m_0 c^2 / \sqrt{1 - v^2/c^2}$. The total energy of the matter is just the sum of all the particle energies, and the flow of this energy through a surface is just the sum of the energy carried by each particle that crosses the surface. We want now to talk only about the energy of the electromagnetic field. So we must write an equation which says that the total *field* energy in a given volume decreases *either* because field energy flows out of the volume *or* because the field loses energy to matter (or gains energy, which is just a negative loss). The field energy inside a volume V is

$$\int_V u \, dV,$$

and its rate of decrease is minus the time derivative of this integral. The flow of field energy out of the volume V is the integral of the normal component of \mathbf{S} over the surface Σ that encloses V ,

$$\int_{\Sigma} \mathbf{S} \cdot \mathbf{n} \, da.$$

So

$$-\frac{\partial}{\partial t} \int_V u \, dV = \int_{\Sigma} \mathbf{S} \cdot \mathbf{n} \, da + (\text{work done on matter inside } V). \quad (27.3)$$

We have seen before that the field does work on each unit volume of matter at the rate $\mathbf{E} \cdot \mathbf{j}$. [The force on a particle is $\mathbf{F} = q(\mathbf{E} + \mathbf{v} \times \mathbf{B})$, and the rate of doing work is $\mathbf{F} \cdot \mathbf{v} = q\mathbf{E} \cdot \mathbf{v}$. If there are N particles per unit volume, the rate of doing work per unit volume is $Nq\mathbf{E} \cdot \mathbf{v}$, but $Nq\mathbf{v} = \mathbf{j}$] So the quantity $\mathbf{E} \cdot \mathbf{j}$ must be equal to the loss of energy per unit time and per unit volume by the field. Equation (27.3) then becomes

$$-\frac{\partial}{\partial t} \int_V u dV = \int_S \mathbf{S} \cdot \mathbf{n} da + \int_V \mathbf{E} \cdot \mathbf{j} dV. \quad (27.4)$$

This is our conservation law for energy in the field. We can convert it into a differential equation like Eq. (27.2) if we can change the second term to a volume integral. That is easy to do with Gauss' theorem. The surface integral of the normal component of \mathbf{S} is the integral of its divergence over the volume inside. So Eq. (27.3) is equivalent to

$$-\int_V \frac{du}{dt} dV = \int_V \nabla \cdot \mathbf{S} dV + \int_V \mathbf{E} \cdot \mathbf{j} dV,$$

where we have put the time derivative of the first term inside the integral. Since this equation is true for any volume, we can take away the integrals and we have the energy equation for the electromagnetic fields:

$$-\frac{\partial u}{\partial t} = \nabla \cdot \mathbf{S} + \mathbf{E} \cdot \mathbf{j}. \quad (27.5)$$

Now this equation doesn't do us a bit of good unless we know what u and \mathbf{S} are. Perhaps we should just tell you what they are in terms of \mathbf{E} and \mathbf{B} , because all we really want is the result. However, we would rather show you the kind of argument that was used by Poynting in 1884 to obtain formulas for \mathbf{S} and u , so you can see where they come from. (You won't, however, need to learn this derivation for our later work.)

27-3 Energy density and energy flow in the electromagnetic field

The idea is to suppose that there is a field energy density u and a flux \mathbf{S} that depend only upon the fields \mathbf{E} and \mathbf{B} . (For example, we know that in electrostatics, at least, the energy density can be written $\frac{1}{2}\epsilon_0 \mathbf{E} \cdot \mathbf{E}$.) Of course, the u and \mathbf{S} might depend on the potentials or something else, but let's see what we can work out. We can try to rewrite the quantity $\mathbf{E} \cdot \mathbf{j}$ in such a way that it becomes the sum of two terms, one that is the time derivative of one quantity and another that is the divergence of a second quantity. The first quantity would then be u and the second would be \mathbf{S} (with suitable signs). Both quantities must be written in terms of the fields only; that is, we want to write our equality as

$$\mathbf{E} \cdot \mathbf{j} = -\frac{\partial u}{\partial t} - \nabla \cdot \mathbf{S}. \quad (27.6)$$

The left-hand side must first be expressed in terms of the fields only. How can we do that? By using Maxwell's equations, of course. From Maxwell's equation for the curl of \mathbf{B} ,

$$\mathbf{j} = \epsilon_0 c^2 \nabla \times \mathbf{B} - \epsilon_0 \frac{\partial \mathbf{E}}{\partial t}.$$

Substituting this in (27.6) we will have only \mathbf{E} 's and \mathbf{B} 's:

$$\mathbf{E} \cdot \mathbf{j} = \epsilon_0 c^2 \mathbf{E} \cdot (\nabla \times \mathbf{B}) - \epsilon_0 \mathbf{E} \cdot \frac{\partial \mathbf{E}}{\partial t}. \quad (27.7)$$

We are already partly finished. The last term is a time derivative—it is $(\partial/\partial t)(\frac{1}{2}\epsilon_0 \mathbf{E} \cdot \mathbf{E})$. So $\frac{1}{2}\epsilon_0 \mathbf{E} \cdot \mathbf{E}$ is at least one part of u . It's the same thing we found in electrostatics. Now, all we have to do is to make the other term into the divergence of something.

Notice that the first term on the right-hand side of (27.7) is the same as

$$(\nabla \times \mathbf{B}) \cdot \mathbf{E}. \quad (27.8)$$

And, as you know from vector algebra, $(\mathbf{a} \times \mathbf{b}) \cdot \mathbf{c}$ is the same as $\mathbf{a} \cdot (\mathbf{b} \times \mathbf{c})$; so our term is also the same as

$$\nabla \cdot (\mathbf{B} \times \mathbf{E}) \quad (27.9)$$

and we have the divergence of “something,” just as we wanted. Only that’s wrong! We warned you before that ∇ is “like” a vector, but not “exactly” the same. The reason it is not is because there is an additional *convention* from calculus: when a derivative operator is in front of a product, it works on everything to the right. In Eq. (27.7), the ∇ operates only on \mathbf{B} , not on \mathbf{E} . But in the form (27.9), the normal convention would say that ∇ operates on both \mathbf{B} and \mathbf{E} . So it’s *not* the same thing. In fact, if we work out the components of $\nabla \cdot (\mathbf{B} \times \mathbf{E})$ we can see that it is equal to $\mathbf{E} \cdot (\nabla \times \mathbf{B})$ *plus* some other terms. It’s like what happens when we take a derivative of a product in algebra. For instance,

$$\frac{d}{dx}(fg) = \frac{df}{dx}g + f\frac{dg}{dx}.$$

Rather than working out all the components of $\nabla \cdot (\mathbf{B} \times \mathbf{E})$, we would like to show you a trick that is very useful for this kind of problem. It is a trick that allows you to use all the rules of vector algebra on expressions with the ∇ operator, without getting into trouble. The trick is to throw out—for a while at least—the rule of the calculus notation about what the derivative operator works on. You see, ordinarily, the order of terms is used for *two* separate purposes. One is for calculus: $f(d/dx)g$ is not the same as $g(d/dx)f$; and the other is for vectors: $\mathbf{a} \times \mathbf{b}$ is different from $\mathbf{b} \times \mathbf{a}$. We can, if we want, choose to abandon momentarily the calculus rule. Instead of saying that a derivative operates on everything to the right, we make a *new* rule that doesn’t depend on the order in which terms are written down. Then we can juggle terms around without worrying.

Here is our new convention. We show, by a subscript, what a differential operator works on; the *order* has no meaning. Suppose we let the operator D stand for $\partial/\partial x$. Then D_f means that only the derivative of the variable quantity f is taken. Then

$$D_f f = \frac{\partial f}{\partial x}.$$

But if we have $D_f fg$, it means

$$D_f fg = \left(\frac{\partial f}{\partial x}\right)g.$$

But notice now that according to our new rule, $f D_f g$ means the same thing. We can write the same thing any which way.

$$D_f fg = g D_f f = f D_f g = fg D_f.$$

You see, the D_f can even come *after* everything. (It’s surprising that such a handy notation is never taught in books on mathematics or physics.)

You may wonder: What if I *want* to write the derivative of fg ? I *want* the derivative of *both* terms. That’s easy, you just say so; you write $D_f(fg) + D_g(fg)$. That is just $g(\partial f/\partial x) + f(\partial g/\partial x)$, which is what you mean in the old notation by $\partial(fg)/\partial x$.

You will see that it is now going to be very easy to work out a new expression for $\nabla \cdot (\mathbf{B} \times \mathbf{E})$. We start by changing to the new notation; we write

$$\nabla \cdot (\mathbf{B} \times \mathbf{E}) = \nabla_B \cdot (\mathbf{B} \times \mathbf{E}) + \nabla_E \cdot (\mathbf{B} \times \mathbf{E}). \quad (27.10)$$

The moment we do that we don’t have to keep the order straight any more. We always know that ∇_E operates on \mathbf{E} only, and ∇_B operates on \mathbf{B} only. In these circumstances, we can use ∇ as though it were an ordinary vector. (Of course,

when we are finished, we will want to return to the “standard” notation that everybody usually uses) So now we can do the various things like interchanging dots and crosses and making other kinds of rearrangements of the terms. For instance, the middle term of Eq. (27.10) can be rewritten as $\mathbf{E} \cdot \nabla_B \times \mathbf{B}$. (You remember that $\mathbf{a} \cdot \mathbf{b} \times \mathbf{c} = \mathbf{b} \cdot \mathbf{c} \times \mathbf{a}$.) And the last term is the same as $\mathbf{B} \cdot \mathbf{E} \times \nabla_E$. It looks freakish, but it is all right. Now if we try to go back to the ordinary convention, we have to arrange that the ∇ operates only on its “own” variable. The first one is already that way, so we can just leave off the subscript. The second one needs some rearranging to put the ∇ in front of the \mathbf{E} , which we can do by reversing the cross product and changing sign:

$$\mathbf{B} \cdot (\mathbf{E} \times \nabla_E) = -\mathbf{B} \cdot (\nabla_E \times \mathbf{E}).$$

Now it is in a conventional order, so we can return to the usual notation. Equation (27.10) is equivalent to

$$\nabla \cdot (\mathbf{B} \times \mathbf{E}) = \mathbf{E} \cdot (\nabla \times \mathbf{B}) - \mathbf{B} \cdot (\nabla \times \mathbf{E}). \quad (27.11)$$

(A quicker way would have been to use components in this special case, but it was worth taking the time to show you the mathematical trick. You probably won’t see it anywhere else, and it is very good for unlocking vector algebra from the rules about the order of terms with derivatives.)

We now return to our energy conservation discussion and use our new result, Eq. (27.11), to transform the $\nabla \times \mathbf{B}$ term of Eq. (27.7). That energy equation becomes

$$\mathbf{E} \cdot \mathbf{j} = \epsilon_0 c^2 \nabla \cdot (\mathbf{B} \times \mathbf{E}) + \epsilon_0 c^2 \mathbf{B} \cdot (\nabla \times \mathbf{E}) - \frac{\partial}{\partial t} \left(\frac{1}{2} \epsilon_0 \mathbf{E} \cdot \mathbf{E} \right) \quad (27.12)$$

Now you see we’re almost finished. We have one term which is a nice derivative with respect to t to use for u and another that is a beautiful divergence to represent \mathbf{S} . Unfortunately, there is the center term left over, which is neither a divergence nor a derivative with respect to t . So we almost made it, but not quite. After some thought, we look back at the differential equations of Maxwell and discover that $\nabla \times \mathbf{E}$ is, fortunately, equal to $-\partial \mathbf{B} / \partial t$, which means that we can turn the extra term into something that is a pure time derivative:

$$\mathbf{B} \cdot (\nabla \times \mathbf{E}) = \mathbf{B} \cdot \left(-\frac{\partial \mathbf{B}}{\partial t} \right) = -\frac{\partial}{\partial t} \left(\frac{\mathbf{B} \cdot \mathbf{B}}{2} \right).$$

Now we have exactly what we want. Our energy equation reads

$$\mathbf{E} \cdot \mathbf{j} = \nabla \cdot (\epsilon_0 c^2 \mathbf{B} \times \mathbf{E}) - \frac{\partial}{\partial t} \left(\frac{\epsilon_0 c^2}{2} \mathbf{B} \cdot \mathbf{B} + \frac{\epsilon_0}{2} \mathbf{E} \cdot \mathbf{E} \right), \quad (27.13)$$

which is exactly like Eq. (27.6), if we make the *definitions*

$$u = \frac{\epsilon_0}{2} \mathbf{E} \cdot \mathbf{E} + \frac{\epsilon_0 c^2}{2} \mathbf{B} \cdot \mathbf{B} \quad (27.14)$$

and

$$\mathbf{S} = \epsilon_0 c^2 \mathbf{E} \times \mathbf{B}. \quad (27.15)$$

(Reversing the cross product makes the signs come out right.)

Our program was successful. We have an expression for the energy density that is the sum of an “electric” energy density and a “magnetic” energy density, whose forms are just like the ones we found in statics *when we worked out the energy in terms of the fields*. Also, we have found a formula for the energy flow vector of the electromagnetic field. This new vector, $\mathbf{S} = \epsilon_0 c^2 \mathbf{E} \times \mathbf{B}$, is called “Poynting’s vector,” after its discoverer. It tells us the rate at which the field energy moves around in space. The energy which flows through a small area da per second is $\mathbf{S} \cdot \mathbf{n} da$, where \mathbf{n} is the unit vector perpendicular to da . (Now that we have our formulas for u and \mathbf{S} , you can forget the derivations if you want.)

27-4 The ambiguity of the field energy

Before we take up some applications of the Poynting formulas [Eqs. (27.14) and (27.15)], we would like to say that we have not really “proved” them. All we did was to find a *possible* “ u ” and a *possible* “ S .” How do we know that by juggling the terms around some more we couldn’t find another formula for “ u ” and another formula for “ S ”? The new S and the new u would be different, but they would still satisfy Eq. (27.6). It’s possible. It can be done, but the forms that have been found always involve various *derivatives* of the field (and always with second-order terms like a second derivative or the square of a first derivative). There are, in fact, an infinite number of different possibilities for u and S , and so far no one has thought of an experimental way to tell which one is right! People have guessed that the simplest one is probably the correct one, but we must say that we do not know for certain what is the actual location in space of the electromagnetic field energy. So we too will take the easy way out and say that the field energy is given by Eq. (27.14). Then the flow vector S must be given by Eq. (27.15).

It is interesting that there seems to be no unique way to resolve the indefiniteness in the location of the field energy. It is sometimes claimed that this problem can be resolved by using the theory of gravitation in the following argument. In the theory of gravity, all energy is the source of gravitational attraction. Therefore the energy density of electricity must be located properly if we are to know in which direction the gravity force acts. As yet, however, no one has done such a delicate experiment that the precise location of the gravitational influence of electromagnetic fields could be determined. That electromagnetic fields alone can be the source of gravitational force is an idea it is hard to do without. It has, in fact, been observed that light is deflected as it passes near the sun—we could say that the sun pulls the light down toward it. Do you not want to allow that the light pulls equally on the sun? Anyway, everyone always accepts the simple expressions we have found for the location of electromagnetic energy and its flow. And although sometimes the results obtained from using them seem strange, nobody has ever found anything wrong with them—that is, no disagreement with experiment. So we will follow the rest of the world—besides, we believe that it is probably perfectly right.

We should make one further remark about the energy formula. In the first place, the energy per unit volume in the field is very simple: It is the electrostatic energy plus the magnetic energy, *if* we write the electrostatic energy in terms of E^2 and the magnetic energy as B^2 . We found two such expressions as *possible* expressions for the energy when we were doing static problems. We also found a number of other formulas for the energy in the electrostatic field, such as $\rho\phi$, which is *equal* to the integral of $\mathbf{E} \cdot \mathbf{E}$ in the electrostatic case. However, in an electrodynamic field the equality failed, and there was no obvious choice as to which was the right one. Now we know which is the right one. Similarly, we have found the formula for the magnetic energy that is correct in general. The right formula for the energy density of *dynamic* fields is Eq. (27.14)

27-5 Examples of energy flow

Our formula for the energy flow vector S is something quite new. We want now to see how it works in some special cases and also to see whether it checks out with anything that we knew before. The first example we will take is light. In a light wave we have an \mathbf{E} vector and a \mathbf{B} vector at right angles to each other and to the direction of the wave propagation. (See Fig. 27-2.) In an electromagnetic wave, the magnitude of \mathbf{B} is equal to $1/c$ times the magnitude of \mathbf{E} , and since they are at right angles,

$$|\mathbf{E} \times \mathbf{B}| = \frac{E^2}{c}$$

Therefore, for light, the flow of energy per unit area per second is

$$S = \epsilon_0 c E^2. \quad (27.16)$$

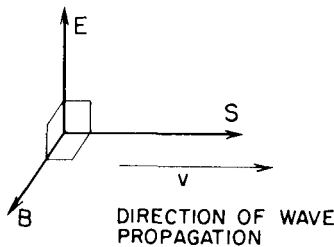


Fig. 27-2. The vectors \mathbf{E} , \mathbf{B} , and \mathbf{S} for a light wave.

For a light wave in which $E = E_0 \cos \omega(t - x/c)$, the average rate of energy flow per unit area, $\langle S \rangle_{av}$ —which is called the “intensity” of the light—is the mean value of the square of the electric field times $\epsilon_0 c$:

$$\text{Intensity} = \langle S \rangle_{av} = \epsilon_0 c \langle E^2 \rangle_{av}. \quad (27.17)$$

Believe it or not, we have already derived this result in Section 31-3 of Vol. I, when we were studying light. We can believe that it is right because it also checks against something else. When we have a light beam, there is an energy density in space given by Eq. (27.14). Using $cB = E$ for a light wave, we get that

$$u = \frac{\epsilon_0}{2} E^2 + \frac{\epsilon_0 c^2}{2} \left(\frac{E^2}{c^2} \right) = \epsilon_0 E^2.$$

But E varies in space, so the average energy density is

$$\langle u \rangle_{av} = \epsilon_0 \langle E^2 \rangle_{av}. \quad (27.18)$$

Now the wave travels at the speed c , so we should think that the energy that goes through a square meter in a second is c times the amount of energy in one cubic meter. So we would say that

$$\langle S \rangle_{av} = \epsilon_0 c \langle E^2 \rangle_{av}.$$

And it's right; it is the same as Eq. (27.17).

Now we take another example. Here is a rather curious one. We look at the energy flow in a capacitor that we are charging slowly. (We don't want frequencies so high that the capacitor is beginning to look like a resonant cavity, but we don't want DC either.) Suppose we use a circular parallel plate capacitor of our usual kind, as shown in Fig. 27-3. There is a nearly uniform electric field inside which is changing with time. At any instant the total electromagnetic energy inside is u times the volume. If the plates have a radius a and a separation h , the total energy between the plates is

$$U = \left(\frac{\epsilon_0}{2} E^2 \right) (\pi a^2 h). \quad (27.19)$$

This energy changes when E changes. When the capacitor is being charged, the volume between the plates is receiving energy at the rate

$$\frac{dU}{dt} = \epsilon_0 \pi a^2 h E \dot{E}. \quad (27.20)$$

So there must be a flow of energy into that volume from somewhere. Of course you know that it must come in on the charging wires—not at all! It can't enter the space between the plates from that direction, because E is perpendicular to the plates; $E \times B$ must be *parallel* to the plates.

You remember, of course, that there is a magnetic field that circles around the axis when the capacitor is charging. We discussed that in Chapter 23. Using the last of Maxwell's equations, we found that the magnetic field at the edge of the capacitor is given by

$$2\pi a c^2 B = \dot{E} \cdot \pi a^2,$$

or

$$B = \frac{a}{2c^2} \dot{E}.$$

Its direction is shown in Fig. 27-3. So there is an energy flow proportional to $E \times B$ that comes in all around the edges, as shown in the figure. The energy isn't actually coming down the wires, but from the space surrounding the capacitor.

Let's check whether or not the total amount of flow through the whole surface between the edges of the plates checks with the rate of change of the energy inside—it had better; we went through all that work proving Eq. (27.15) to make sure,

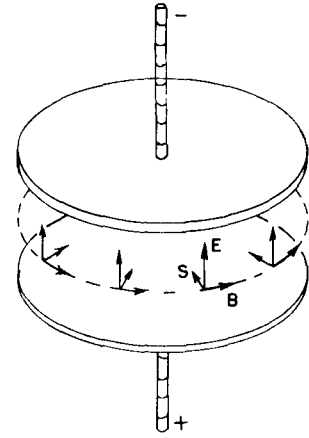


Fig. 27-3. Near a charging capacitor, the Poynting vector S points inward toward the axis.

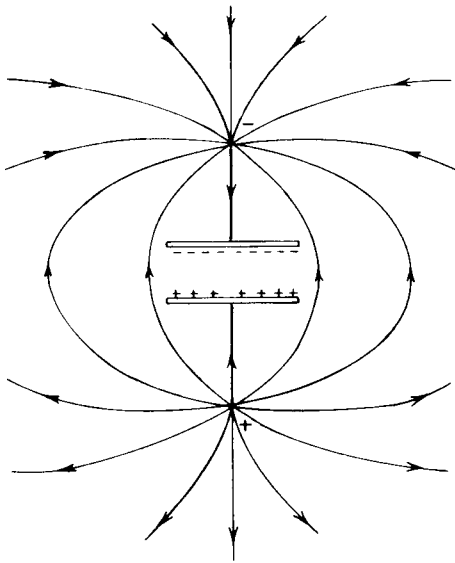


Fig. 27-4. The fields outside a capacitor when it is being charged by bringing two charges from a large distance.

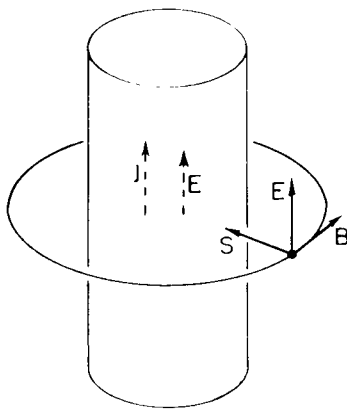


Fig. 27-5 The Poynting vector \mathbf{S} near a wire carrying a current.

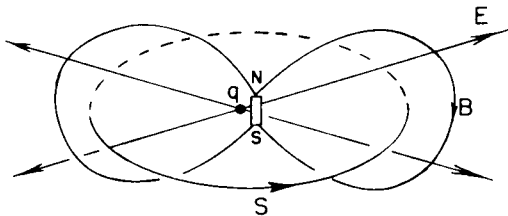


Fig. 27-6. A charge and a magnet produce a Poynting vector that circulates in closed loops.

but let's see. The area of the surface is $2\pi ah$, and $\mathbf{S} = \epsilon_0 c^2 \mathbf{E} \times \mathbf{B}$ is in magnitude

$$\epsilon_0 c^2 E \left(\frac{a}{2c^2} \dot{E} \right),$$

so the total flux of energy is

$$\pi a^2 h \epsilon_0 E \dot{E}.$$

It does check with Eq. (27.20). But it tells us a peculiar thing: that when we are charging a capacitor, the energy is not coming down the wires; it is coming in through the edges of the gap. That's what this theory says!

How can that be? That's *not* an easy question, but here is one way of thinking about it. Suppose that we had some charges above and below the capacitor and far away. When the charges are far away, there is a weak but enormously spread-out field that surrounds the capacitor. (See Fig. 27-4.) Then, as the charges come together, the field gets stronger nearer to the capacitor. So the field energy which is way out moves toward the capacitor and eventually ends up between the plates.

As another example, we ask what happens in a piece of resistance wire when it is carrying a current. Since the wire has resistance, there is an electric field along it, driving the current. Because there is a potential drop along the wire, there is also an electric field just outside the wire, parallel to the surface. (See Fig. 27-5.) There is, in addition, a magnetic field which goes around the wire because of the current. The \mathbf{E} and \mathbf{B} are at right angles; therefore there is a Poynting vector directed radially inward, as shown in the figure. There is a flow of energy into the wire all around. It is, of course, equal to the energy being lost in the wire in the form of heat. So our "crazy" theory says that the electrons are getting their energy to generate heat because of the energy flowing into the wire from the field outside. Intuition would seem to tell us that the electrons get their energy from being pushed along the wire, so the energy should be flowing down (or up) along the wire. But the theory says that the electrons are really being pushed by an electric field, which has come from some charges very far away, and that the electrons get their energy for generating heat from these fields. The energy somehow flows from the distant charges into a wide area of space and then inward to the wire.

Finally, in order to really convince you that this theory is obviously nuts, we will take one more example—an example in which an electric charge and a magnet are *at rest* near each other—both sitting quite still. Suppose we take the example of a point charge sitting near the center of a bar magnet, as shown in Fig. 27-6. Everything is at rest, so the energy is not changing with time. Also, \mathbf{E} and \mathbf{B} are quite static. But the Poynting vector says that there is a flow of energy, because there is an $\mathbf{E} \times \mathbf{B}$ that is not zero. If you look at the energy flow, you find that it just circulates around and around. There isn't any change in the energy anywhere—everything which flows into one volume flows out again. It is like incompressible water flowing around. So there is a circulation of energy in this so-called static condition. How absurd it gets!

Perhaps it isn't so terribly puzzling, though, when you remember that what we called a "static" magnet is really a circulating permanent current. In a permanent magnet the electrons are spinning permanently inside. So maybe a circulation of the energy outside isn't so queer after all.

You no doubt begin to get the impression that the Poynting theory at least partially violates your intuition as to where energy is located in an electromagnetic field. You might believe that you must revamp all your intuitions, and, therefore have a lot of things to study here. But it seems really not necessary. You don't need to feel that you will be in great trouble if you forget once in a while that the energy in a wire is flowing into the wire from the outside, rather than along the wire. It seems to be only rarely of value, when using the idea of energy conservation, to notice in detail what path the energy is taking. The circulation of energy around a magnet and a charge seems, in most circumstances, to be quite unimportant. It is not a vital detail, but it is clear that our ordinary intuitions are quite wrong.

27-6 Field momentum

Next we would like to talk about the *momentum* in the electromagnetic field. Just as the field has energy, it will have a certain momentum per unit volume. Let us call that momentum density \mathbf{g} . Of course, momentum has various possible directions, so that \mathbf{g} must be a vector. Let's talk about one component at a time; first, we take the x -component. Since each component of momentum is conserved we should be able to write down a law that looks something like this:

$$-\frac{\partial}{\partial t} \left(\text{momentum of matter} \right)_x = \frac{\partial g_x}{\partial t} + \left(\text{momentum outflow} \right)_x.$$

The left side is easy. The rate-of-change of the momentum of matter is just the force on it. For a particle, it is $\mathbf{F} = q(\mathbf{E} + \mathbf{v} \times \mathbf{B})$; for a distribution of charges, the force per unit volume is $(\rho\mathbf{E} + \mathbf{j} \times \mathbf{B})$. The "momentum outflow" term, however, is strange. It cannot be the divergence of a vector because it is not a scalar; it is, rather, an x -component of some vector. Anyway, it should probably look something like

$$\frac{\partial a}{\partial x} + \frac{\partial b}{\partial y} + \frac{\partial c}{\partial z},$$

because the x -momentum could be flowing in any one of the three directions. In any case, whatever a , b , and c are, the combination is supposed to equal the outflow of the x -momentum.

Now the game would be to write $\rho\mathbf{E} + \mathbf{j} \times \mathbf{B}$ in terms only of \mathbf{E} and \mathbf{B} —eliminating ρ and \mathbf{j} by using Maxwell's equations—and then to juggle terms and make substitutions to get it into a form that looks like

$$\frac{\partial g_x}{\partial t} + \frac{\partial a}{\partial x} + \frac{\partial b}{\partial y} + \frac{\partial c}{\partial z}.$$

Then, by identifying terms, we would have expressions for g_x , a , b , and c . It's a lot of work, and we are not going to do it. Instead, we are only going to find an expression for \mathbf{g} , the momentum density—and by a different route.

There is an important theorem in mechanics which is this: whenever there is a flow of energy in any circumstance at all (field energy or any other kind of energy), the energy flowing through a unit area per unit time, when multiplied by $1/c^2$, is equal to the momentum per unit volume in the space. In the special case of electrodynamics, this theorem gives the result that \mathbf{g} is $1/c^2$ times the Poynting vector

$$\mathbf{g} = \frac{1}{c^2} \mathbf{S}. \quad (27.21)$$

So the Poynting vector gives not only energy flow but, if you divide by c^2 , also the momentum density. The same result would come out of the other analysis we suggested, but it is more interesting to notice this more general result. We will now give a number of interesting examples and arguments to convince you that the general theorem is true.

First example: Suppose that we have a lot of particles in a box—let's say N per cubic meter—and that they are moving along with some velocity \mathbf{v} . Now let's consider an imaginary plane surface perpendicular to \mathbf{v} . The energy flow through a unit area of this surface per second is equal to Nv , the number which flow through the surface per second, times the energy carried by each one. The energy in each particle is $m_0c^2/\sqrt{1 - v^2/c^2}$. So the energy flow per second is

$$Nv \frac{m_0c^2}{\sqrt{1 - v^2/c^2}}.$$

But the momentum of each particle is $m_0v/\sqrt{1 - v^2/c^2}$, so the *density* of momentum is

$$N \frac{m_0v}{\sqrt{1 - v^2/c^2}},$$

which is just $1/c^2$ times the energy flow—as the theorem says. So the theorem is true for a bunch of particles.

It is also true for light. When we studied light in Volume I, we saw that when the energy is absorbed from a light beam, a certain amount of momentum is delivered to the absorber. We have, in fact, shown in Chapter 36 of Vol. I that the momentum is $1/c$ times the energy absorbed [Eq. (36.24) of Vol. I]. If we let U_0 be the energy arriving at a unit area per second, then the momentum arriving at a unit area per second is U_0/c . But the momentum is travelling at the speed c , so its *density* in front of the absorber must be U_0/c^2 . So again the theorem is right.

Finally we will give an argument due to Einstein which demonstrates the same thing once more. Suppose that we have a railroad car on wheels (assumed frictionless) with a certain big mass M . At one end there is a device which will shoot out some particles or light (or anything, it doesn't make any difference what it is), which are then stopped at the opposite end of the car. There was some energy originally at one end—say the energy U indicated in Fig. 27-7(a)—and then later it is at the opposite end, as shown in Fig. 27-7(c). The energy U has been displaced the distance L , the length of the car. Now the energy U has the mass U/c^2 , so if the car stayed still, the center of gravity of the car would be moved. Einstein didn't like the idea that the center of gravity of an object could be moved by fooling around only on the inside, so he assumed that it is impossible to move the center of gravity by doing anything inside. But if that is the case, when we moved the energy U from one end to the other, the whole car must have recoiled some distance x , as shown in part (c) of the figure. You can see, in fact, that the total mass of the car, times x , must equal the mass of the energy moved, U/c^2 times L (assuming that U/c^2 is much less than M):

$$Mx = \frac{U}{c^2} L. \quad (27.22)$$

Let's now look at the special case of the energy being carried by a light flash. (The argument would work as well for particles, but we will follow Einstein, who was interested in the problem of light.) What causes the car to be moved? Einstein argued as follows: When the light is emitted there must be a recoil, some unknown recoil with momentum p . It is this recoil which makes the car roll backward. The recoil velocity v of the car will be this momentum divided by the mass of the car:

$$v = \frac{p}{M}.$$

The car moves with this velocity until the light energy U gets to the opposite end. Then, when it hits, it gives back its momentum and stops the car. If x is small, then the time the car moves is nearly equal to L/c ; so we have that

$$x = vt = v \frac{L}{c} = \frac{p}{M} \frac{L}{c}.$$

Putting this x in Eq. (27.22), we get that

$$p = \frac{U}{c}.$$

Again we have the relation of energy and momentum for light. Dividing by c to get the momentum density $g = p/c$, we get once more that

$$g = \frac{U}{c^2}. \quad (27.23)$$

You may well wonder: What is so important about the center-of-gravity theorem? Maybe it is wrong. Perhaps, but then we would also lose the conservation of angular momentum. Suppose that our boxcar is moving along a track at some speed v and that we shoot some light energy from the *top* to the *bottom* of the car—say, from A to B in Fig. 27-8. Now we look at the angular momentum of the system about the point P . Before the energy U leaves A , it has the mass

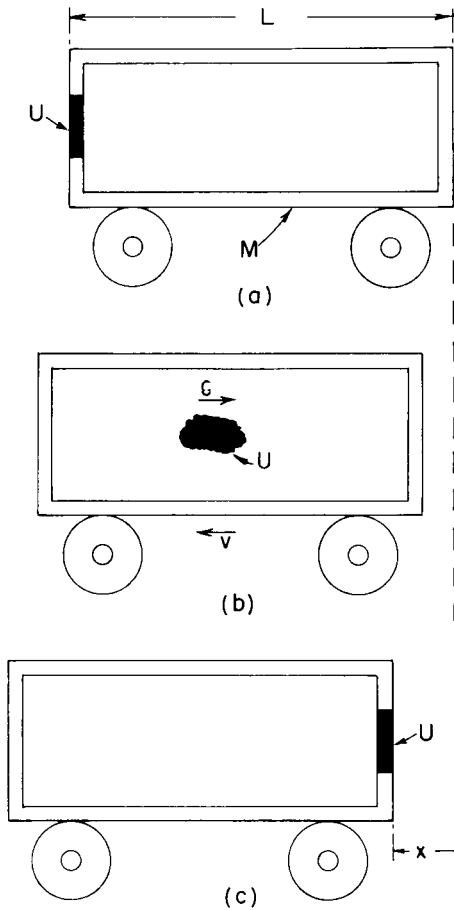


Fig. 27-7. The energy U in motion at the speed c carries the momentum U/c .

$m = U^2/c$ and the velocity v , so it has the angular momentum mvr_a . When it arrives at B , it has the same mass and, if the *linear* momentum of the whole boxcar is not to change, it must still have the velocity v . Its angular momentum about P is then mvr_B . The angular momentum will be changed *unless* the right recoil momentum was given to the car when the light was emitted—that is, unless the light carries the momentum U/c . It turns out that the angular momentum conservation and the theorem of center-of-gravity are closely related in the relativity theory. So the conservation of angular momentum would also be destroyed if our theorem were not true. At any rate, it does turn out to be a true general law, and in the case of electrodynamics we can use it to get the momentum in the field.

We will mention two further examples of momentum in the electromagnetic field. We pointed out in Section 26-2 the failure of the law of action and reaction when two charged particles were moving on orthogonal trajectories. The forces on the two particles don't balance out, so the action and reaction are not equal, therefore the net momentum of the matter must be changing. It is not conserved. But the momentum in the field is also changing in such a situation. If you work out the amount of momentum given by the Poynting vector, it is not constant. However, the change of the particle momenta is just made up by the field momentum, so the total momentum of particles plus field is conserved.

Finally, another example is the situation with the magnet and the charge, shown in Fig. 27-6. We were unhappy to find that energy was flowing around in circles, but now, since we know that energy flow and momentum are proportional, we know also that there is momentum circulating in the space. But a *circulating* momentum means that there is *angular* momentum. So there is *angular* momentum in the field. Do you remember the paradox we described in Section 17-4 about a solenoid and some charges mounted on a disc? It seemed that when the current turned off, the whole disc should start to turn. The puzzle was: Where did the angular momentum come from? The answer is that if you have a magnetic field and some charges, there will be some angular momentum in the field. It must have been put there when the field was built up. When the field is turned off, the angular momentum is given back. So the disc in the paradox *would* start rotating. This mystic circulating flow of energy, which at first seemed so ridiculous, is absolutely necessary. There is really a momentum flow. It is needed to maintain the conservation of angular momentum in the whole world.

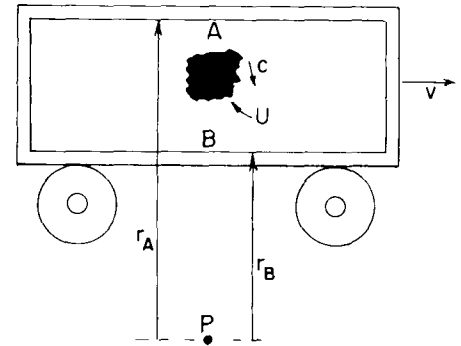


Fig. 27-8. The energy U must carry the momentum U/c if the angular momentum about P is to be conserved.

Electromagnetic Mass

28-1 The field energy of a point charge

In bringing together relativity and Maxwell's equations, we have finished our main work on the theory of electromagnetism. There are, of course, some details we have skipped over and one large area that we will be concerned with in the future—the interaction of electromagnetic fields with matter. But we want to stop for a moment to show you that this tremendous edifice, which is such a beautiful success in explaining so many phenomena, ultimately falls on its face. When you follow any of our physics too far, you find that it always gets into some kind of trouble. Now we want to discuss a serious trouble—the failure of the classical electromagnetic theory. You can appreciate that there is a failure of all classical physics because of the quantum-mechanical effects. Classical mechanics is a mathematically consistent theory; it just doesn't agree with experience. It is interesting, though, that the classical theory of electromagnetism is an unsatisfactory theory all by itself. There are difficulties associated with the *ideas* of Maxwell's theory which are not solved by and not directly associated with quantum mechanics. You may say, "Perhaps there's no use worrying about these difficulties. Since the quantum mechanics is going to change the laws of electrodynamics, we should wait to see what difficulties there are after the modification." However, when electromagnetism is joined to quantum mechanics, the difficulties remain. So it will not be a waste of our time now to look at what these difficulties are. Also, they are of great historical importance. Furthermore, you may get some feeling of accomplishment from being able to go far enough with the theory to see everything—including all of its troubles.

The difficulty we speak of is associated with the concepts of electromagnetic momentum and energy, when applied to the electron or any charged particle. The concepts of simple charged particles and the electromagnetic field are in some way inconsistent. To describe the difficulty, we begin by doing some exercises with our energy and momentum concepts.

First, we compute the energy of a charged particle. Suppose we take a simple model of an electron in which all of its charge q is uniformly distributed on the surface of a sphere of radius a , which we may take to be zero for the special case of a point charge. Now let's calculate the energy in the electromagnetic field. If the charge is standing still, there is no magnetic field, and the energy per unit volume is proportional to the square of the electric field. The magnitude of the electric field is $q/4\pi\epsilon_0 r^2$, and the energy density is

$$u = \frac{\epsilon_0}{2} E^2 = \frac{q^2}{32\pi^2\epsilon_0 r^4}.$$

To get the total energy, we must integrate this density over all space. Using the volume element $4\pi r^2 dr$, the total energy, which we will call U_{elec} , is

$$U_{\text{elec}} = \int \frac{q^2}{8\pi\epsilon_0 r^2} dr.$$

This is readily integrated. The lower limit is a , and the upper limit is ∞ , so

$$U_{\text{elec}} = \frac{1}{2} \frac{q^2}{4\pi\epsilon_0} \frac{1}{a}. \quad (28.1)$$

28-1 The field energy of a point charge

28-2 The field momentum of a moving charge

28-3 Electromagnetic mass

28-4 The force of an electron on itself

28-5 Attempts to modify the Maxwell theory

28-6 The nuclear force field

If we use the electronic charge q_e for q and the symbol e^2 for $q_e^2/4\pi\epsilon_0$, then

$$U_{\text{elec}} = \frac{1}{2} \frac{e^2}{a}. \quad (28.2)$$

It is all fine until we set a equal to zero for a point charge—there's the great difficulty. Because the energy of the field varies inversely as the fourth power of the distance from the center, its volume integral is infinite. There is an infinite amount of energy in the field surrounding a point charge.

What's wrong with an infinite energy? If the energy can't get out, but must stay there forever, is there any real difficulty with an infinite energy? Of course, a quantity that comes out infinite may be annoying, but what really matters is only whether there are any *observable* physical effects. To answer that question, we must turn to something else besides the energy. Suppose we ask how the energy *changes* when we *move* the charge. Then, if the *changes* are infinite, we will be in trouble.

28-2 The field momentum of a moving charge

Suppose an electron is moving at a uniform velocity through space, assuming for a moment that the velocity is low compared with the speed of light. Associated with this moving electron there is a momentum—even if the electron had no mass before it was charged—because of the momentum in the electromagnetic field. We can show that the field momentum is in the direction of the velocity v of the charge and is, for small velocities, proportional to v . For a point P at the distance r from the center of the charge and at the angle θ with respect to the line of motion (see Fig. 28-1) the electric field is radial and, as we have seen, the magnetic field is $v \times E/c^2$. The momentum density, Eq. (27.21), is

$$g = \epsilon_0 E \times B.$$

It is directed obliquely toward the line of motion, as shown in the figure, and has the magnitude

$$g = \frac{\epsilon_0 v}{c^2} E^2 \sin \theta.$$

The fields are symmetric about the line of motion, so when we integrate over space, the transverse components will sum to zero, giving a resultant momentum parallel to v . The component of g in this direction is $g \sin \theta$, which we must integrate over all space. We take as our volume element a ring with its plane perpendicular to v , as shown in Fig. 28-2. Its volume is $2\pi r^2 \sin \theta d\theta dr$. The total momentum is then

$$p = \int \frac{\epsilon_0 v}{c^2} E^2 \sin^2 \theta 2\pi r^2 \sin \theta d\theta dr.$$

Since E is independent of θ (for $v \ll c$), we can immediately integrate over θ ; the integral is

$$\int \sin^3 \theta d\theta = - \int (1 - \cos^2 \theta) d(\cos \theta) = -\cos \theta + \frac{\cos^3 \theta}{3}.$$

The limits of θ are 0 and π , so the θ -integral gives merely a factor of $4/3$, and

$$p = \frac{8\pi}{3} \frac{\epsilon_0 v}{c^2} \int E^2 r^2 dr.$$

The integral (for $v \ll c$) is the one we have just evaluated to find the energy; it is $q^2/16\pi^2\epsilon_0 a$, and

$$p = \frac{2}{3} \frac{q^2}{4\pi\epsilon_0} \frac{v}{ac^2},$$

or

$$p = \frac{2}{3} \frac{e^2}{ac^2} v. \quad (28.3)$$

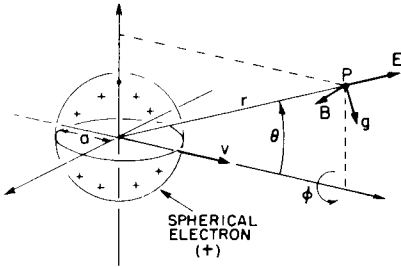


Fig. 28-1. The fields E and B and the momentum density g for a positive electron. For a negative electron, E and B are reversed but g is not.

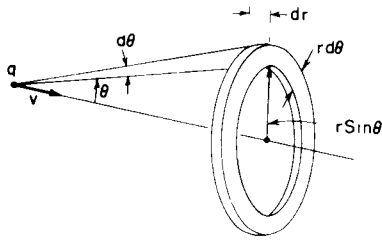


Fig. 28-2. The volume element $2\pi r^2 \sin \theta d\theta dr$ used for calculating the field momentum.

The momentum in the field—the electromagnetic momentum—is proportional to v . It is just what we should have for a particle with the mass equal to the coefficient of v . We can, therefore, call this coefficient the *electromagnetic mass*, m_{elec} , and write it as

$$m_{\text{elec}} = \frac{2}{3} \frac{e^2}{ac^2}. \quad (28.4)$$

28-3 Electromagnetic mass

Where does the mass come from? In our laws of mechanics we have supposed that every object “carries” a thing we call the mass—which also means that it “carries” a momentum proportional to its velocity. Now we discover that it is understandable that a charged particle carries a momentum proportional to its velocity. It might, in fact, be that the mass is just the effect of electrodynamics. The origin of mass has until now been unexplained. We have at last in the theory of electrodynamics a grand opportunity to understand something that we never understood before. It comes out of the blue—or rather, from Maxwell and Poynting—that any charged particle will have a momentum proportional to its velocity just from electromagnetic influences.

Let's be conservative and say, for a moment, that there are two kinds of mass—that the total momentum of an object could be the sum of a mechanical momentum and the electromagnetic momentum. The mechanical momentum is the “mechanical” mass, m_{mech} , times v . In experiments where we measure the mass of a particle by seeing how much momentum it has, or how it swings around in an orbit, we are measuring the total mass. We say generally that the momentum is the total mass ($m_{\text{mech}} + m_{\text{elec}}$) times the velocity. So the observed mass can consist of two pieces (or possibly more if we include other fields): a mechanical piece plus an electromagnetic piece. We know that there is definitely an electromagnetic piece, and we have a formula for it. And there is the thrilling possibility that the mechanical piece is not there at all—that the mass is all electromagnetic.

Let's see what size the electron must have if there is to be no mechanical mass. We can find out by setting the electromagnetic mass of Eq. (28.4) equal to the observed mass m_e of an electron. We find

$$a = \frac{2}{3} \frac{e^2}{m_e c^2}. \quad (28.5)$$

The quantity

$$r_0 = \frac{e^2}{m_e c^2} \quad (28.6)$$

is called the “classical electron radius”; it has the numerical value 2.82×10^{-13} cm, about one one-hundred-thousandth of the diameter of an atom.

Why is r_0 called the electron radius, rather than our a ? Because we could equally well do the same calculation with other assumed distributions of charges—the charge might be spread uniformly through the volume of a sphere or it might be smeared out like a fuzzy ball. For any particular assumption the factor $2/3$ would change to some other fraction. For instance, for a charge uniformly distributed throughout the volume of a sphere, the $2/3$ gets replaced by $4/5$. Rather than to argue over which distribution is correct, it was decided to define r_0 as the “nominal” radius. Then different theories could supply their pet coefficients.

Let's pursue our electromagnetic theory of mass. Our calculation was for $v \ll c$; what happens if we go to high velocities? Early attempts led to a certain amount of confusion, but Lorentz realized that the charged sphere would contract into an ellipsoid at high velocities and that the fields would change in accordance with the formulas (26.6) and (26.7) we derived for the relativistic case in Chapter 26. If you carry through the integrals for \mathbf{p} in that case, you find that for an arbitrary velocity v , the momentum is altered by the factor $1/\sqrt{1 - v^2/c^2}$:

$$\mathbf{p} = \frac{2}{3} \frac{e^2}{ac^2} \frac{\mathbf{v}}{\sqrt{1 - v^2/c^2}}. \quad (28.7)$$

In other words, the electromagnetic mass rises with velocity inversely as $\sqrt{1 - v^2/c^2}$ —a discovery that was made before the theory of relativity.

Early experiments were proposed to measure the changes with velocity in the observed mass of a particle in order to determine how much of the mass was mechanical and how much was electrical. It was believed at the time that the electrical part *would* vary with velocity, whereas the mechanical part would *not*. But while the experiments were being done, the theorists were also at work. Soon the theory of relativity was developed, which proposed that no matter what the origin of the mass, it *all* should vary as $m_0/\sqrt{1 - v^2/c^2}$. Equation (28.7) was the beginning of the theory that mass depended on velocity.

Let's now go back to our calculation of the energy in the field, which led to Eq. (28.2). According to the theory of relativity, the energy U will have the mass U/c^2 ; Eq. (28.2) then says that the field of the electron should have the mass

$$m'_{\text{elec}} = \frac{U_{\text{elec}}}{c^2} = \frac{1}{2} \frac{e^2}{ac^2}, \quad (28.8)$$

which is not the same as the electromagnetic mass, m_{elec} , of Eq. (28.4). In fact, if we just combine Eqs. (28.2) and (28.4), we would write

$$U_{\text{elec}} = \frac{3}{4} m_{\text{elec}} c^2.$$

This formula was discovered before relativity, and when Einstein and others began to realize that it must always be that $U = mc^2$, there was great confusion.

28-4 The force of an electron on itself

The discrepancy between the two formulas for the electromagnetic mass is especially annoying, because we have carefully proved that the theory of electrodynamics is consistent with the principle of relativity. Yet the theory of relativity implies without question that the momentum must be the same as the energy times v/c^2 . So we are in some kind of trouble; we must have made a mistake. We did not make an algebraic mistake in our calculations, but we have left something out.

In deriving our equations for energy and momentum, we assumed the conservation laws. We assumed that *all* forces were taken into account and that any work done and any momentum carried by other “nonelectrical” machinery was included. Now if we have a sphere of charge, the electrical forces are all repulsive and an electron would tend to fly apart. Because the system has unbalanced forces, we can get all kinds of errors in the laws relating energy and momentum. To get a *consistent* picture, we must imagine that something holds the electron together. The charges must be *held* to the sphere by some kind of rubber bands—something that keeps the charges from flying off. It was first pointed out by Poincaré that the rubber bands—or whatever it is that holds the electron together—must be included in the energy and momentum calculations. For this reason the extra nonelectrical forces are also known by the more elegant name “the Poincaré stresses.” If the extra forces are included in the calculations, the masses obtained in two ways are changed (in a way that depends on the detailed assumptions). And the results are consistent with relativity; i.e., the mass that comes out from the momentum calculation is the same as the one that comes from the energy calculation. However, both of them contain *two* contributions: an electromagnetic mass and contribution from the Poincaré stresses. Only when the two are added together do we get a consistent theory.

It is therefore impossible to get all the mass to be electromagnetic in the way we hoped. It is not a legal theory if we have nothing but electrodynamics. Something else has to be added. Whatever you call them—“rubber bands,” or “Poincaré stresses,” or something else—there have to be other forces in nature to make a consistent theory of this kind.

Clearly, as soon as we have to put forces on the inside of the electron, the beauty of the whole idea begins to disappear. Things get very complicated. You would want to ask: How strong are the stresses? How does the electron shake? Does it oscillate? What are all its internal properties? And so on. It might be possible that an electron does have some complicated internal properties. If we made a theory of the electron along these lines, it would predict odd properties, like modes of oscillation, which haven't apparently been observed. We say "apparently" because we observe a lot of things in nature that still do not make sense. We may someday find out that one of the things we don't understand today (for example, the muon) can, in fact, be explained as an oscillation of the Poincaré stresses. It doesn't seem likely, but no one can say for sure. There are so many things about fundamental particles that we still don't understand. Anyway, the complex structure implied by this theory is undesirable, and the attempt to explain all mass in terms of electromagnetism—at least in the way we have described—has led to a blind alley.

We would like to think a little more about why we say we have a mass when the momentum in the field is proportional to the velocity. Easy! The mass is the coefficient between momentum and velocity. But we can look at the mass in another way: a particle has mass if you have to exert a force in order to accelerate it. So it may help our understanding if we look a little more closely at where the forces come from. How do we know that there has to be a force? Because we have proved the law of the conservation of momentum for the fields. If we have a charged particle and push on it for awhile, there will be some momentum in the electromagnetic field. Momentum must have been poured into the field somehow. Therefore there must have been a force pushing on the electron in order to get it going—a force in addition to that required by its mechanical inertia, a force due to its electromagnetic interaction. And there must be a corresponding force back on the "pusher." But where does that force come from?

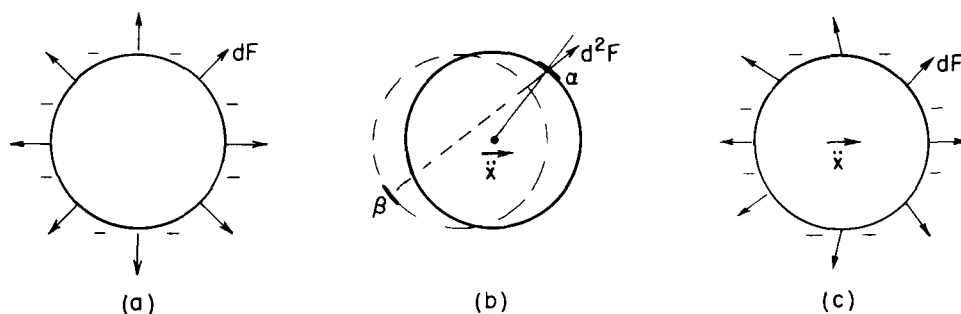


Fig. 28-3. The self-force on an accelerating electron is not zero because of the retardation. (By dF we mean the force on a surface element da ; by d^2F we mean the force on the surface element da_α from the charge on the surface element da_β .)

The picture is something like this. We can think of the electron as a charged sphere. When it is at rest, each piece of charge repels electrically each other piece, but the forces all balance in pairs, so that there is no *net* force. [See Fig. 28-3(a).] However, when the electron is being accelerated, the forces will no longer be in balance because of the fact that the electromagnetic influences take time to go from one piece to another. For instance, the force on the piece α in Fig. 28-3(b) from a piece β on the opposite side depends on the position of β at an earlier time, as shown. Both the magnitude and direction of the force depend on the motion of the charge. If the charge is accelerating, the forces on various parts of the electron might be as shown in Fig. 28-3(c). When all these forces are added up, they don't cancel out. They would cancel for a uniform velocity, even though it looks at first glance as though the retardation would give an unbalanced force even for a uniform velocity. But it turns out that there is no net force unless the electron is being accelerated. With acceleration, if we look at the forces between

the various parts of the electron, action and reaction are not exactly equal, and the electron exerts a force *on itself* that tries to hold back the acceleration. It holds itself back by its own bootstraps.

It is possible, but difficult, to calculate this self-reaction force; however, we don't want to go into such an elaborate calculation here. We will tell you what the result is for the special case of relatively uncomplicated motion in one dimension, say x . Then, the self-force can be written in a series. The first term in the series depends on the acceleration \ddot{x} , the next term is proportional to \dot{x} , and so on.* The result is

$$F = \alpha \frac{e^2}{ac^2} \dot{x} - \frac{2}{3} \frac{e^2}{c^3} \ddot{x} + \gamma \frac{e^2 a}{c^4} \ddot{\ddot{x}} + \cdots \quad (28.9)$$

where α and γ are numerical coefficients of the order of 1. The coefficient α of the \ddot{x} term depends on what charge distribution is assumed; if the charge is distributed uniformly on a sphere, then $\alpha = 2/3$. So there is a term, proportional to the acceleration, which varies inversely as the radius a of the electron and agrees exactly with the value we got in Eq. (28.4) for m_{elec} . If the charge distribution is chosen to be different, so that α is changed, the fraction $2/3$ in Eq. (28.4) would be changed in the same way. The term in \ddot{x} is *independent* of the assumed radius a , and also of the assumed distribution of the charge; its coefficient is *always* $2/3$. The next term is proportional to the radius a , and its coefficient γ depends on the charge distribution. You will notice that if we let the electron radius a go to zero, the last term (and all higher terms) will go to zero; the second term remains constant, but the first term—the electromagnetic mass—goes to infinity. And we can see that the infinity arises because of the force of one part of the electron on another—because we have allowed what is perhaps a silly thing, the possibility of the “point” electron acting on itself.

28-5 Attempts to modify the Maxwell theory

We would like now to discuss how it might be possible to modify Maxwell's theory of electrodynamics so that the idea of an electron as a simple point charge could be maintained. Many attempts have been made, and some of the theories were even able to arrange things so that all the electron mass was electromagnetic. But all of these theories have died. It is still interesting to discuss some of the possibilities that have been suggested—to see the struggles of the human mind.

We started out our theory of electricity by talking about the interaction of one charge with another. Then we made up a theory of these interacting charges and ended up with a field theory. We believe it so much that we allow it to tell us about the force of one part of an electron on another. Perhaps the entire difficulty is that electrons do not act on themselves; perhaps we are making too great an extrapolation from the interaction of separate electrons to the idea that an electron interacts with itself. Therefore some theories have been proposed in which the possibility that an electron acts on itself is ruled out. Then there is no longer the infinity due to the self-action. Also, there is no longer any electromagnetic mass associated with the particle; all the mass is back to being mechanical, but there are new difficulties in the theory.

We must say immediately that such theories require a modification of the idea of the electromagnetic field. You remember we said at the start that the force on a particle at any point was determined by just two quantities— \mathbf{E} and \mathbf{B} . If we abandon the “self-force” this can no longer be true, because if there is an electron in a certain place, the force isn't given by the total \mathbf{E} and \mathbf{B} , but by only those parts due to *other* charges. So we have to keep track always of how much of \mathbf{E} and \mathbf{B} is due to the charge on which you are calculating the force and how much is due to the other charges. This makes the theory much more elaborate, but it gets rid of the difficulty of the infinity.

* We are using the notation: $x = dx/dt$, $\dot{x} = d^2x/dt^2$, $\ddot{x} = d^3x/dt^3$, etc.

So we can, *if we want to*, say that there is no such thing as the electron acting upon itself, and throw away the whole set of forces in Eq. (28.9). However, we have then thrown away the baby with the bath! Because the second term in Eq. (28.9), the term in \ddot{x} , is needed. That force does something very definite. If you throw it away, you're in trouble again. When we accelerate a charge, it radiates electromagnetic waves, so it loses energy. Therefore, to accelerate a charge, we must require more force than is required to accelerate a neutral object of the same mass; otherwise energy wouldn't be conserved. The rate at which we do work on an accelerating charge must be equal to the rate of loss of energy per second by radiation. We have talked about this effect before—it is called the radiation resistance. We still have to answer the question: Where does the extra force, against which we must do this work, come from? When a big antenna is radiating, the forces come from the influence of one part of the antenna current on another. For a single accelerating electron radiating into otherwise empty space, there would seem to be only one place the force could come from—the action of one part of the electron on another part.

We found back in Chapter 32 of Vol. I that an oscillating charge radiates energy at the rate

$$\frac{dW}{dt} = \frac{2}{3} \frac{e^2 (\ddot{x})^2}{c^3}. \quad (28.10)$$

Let's see what we get for the rate of doing work *on* an electron against the bootstrap force of Eq. (28.9). The rate of work is the force times the velocity, or $F \dot{x}$:

$$\frac{dW}{dt} = \alpha \frac{e^2}{ac^2} \ddot{x} \dot{x} - \frac{2}{3} \frac{e^2}{c^3} \ddot{x} \dot{x} + \dots \quad (28.11)$$

The first term is proportional to $d\dot{x}^2/dt$, and therefore just corresponds to the rate of change of the kinetic energy $\frac{1}{2} m v^2$ associated with the electromagnetic mass. The second term should correspond to the radiated power in Eq. (28.10). But it is different. The discrepancy comes from the fact that the term in Eq. (28.11) is generally true, whereas Eq. (28.10) is right only for an *oscillating* charge. We can show that the two are equivalent if the motion of the charge is periodic. To do that, we rewrite the second term of Eq. (28.11) as

$$- \frac{2}{3} \frac{e^2}{c^3} \frac{d}{dt} (\dot{x} \ddot{x}) + \frac{2}{3} \frac{e^2}{c^3} (\ddot{x})^2,$$

which is just an algebraic transformation. If the motion of the electron is periodic, the quantity $\dot{x} \ddot{x}$ returns periodically to the same value, so that if we take the *average* of its time derivative, we get zero. The second term, however, is always positive (it's a square), so its average is also positive. This term gives the net work done and is just equal to Eq. (28.10).

The term in \ddot{x} of the bootstrap force is required in order to have energy conservation in radiating systems, and we can't throw it away. It was, in fact, one of the triumphs of Lorentz to show that there is such a force and that it comes from the action of the electron on itself. We must believe in the idea of the action of the electron on itself, and we *need* the term in \ddot{x} . The problem is how we can get that term without getting the first term in Eq. (28.9), which gives all the trouble. We don't know how. You see that the classical electron theory has pushed itself into a tight corner.

There have been several other attempts to modify the laws in order to straighten the thing out. One way, proposed by Born and Infeld, is to change the Maxwell equations in a complicated way so that they are no longer linear. Then the electromagnetic energy and momentum can be made to come out finite. But the laws they suggest predict phenomena which have never been observed. Their theory also suffers from another difficulty we will come to later, which is common to all the attempts to avoid the troubles we have described.

The following peculiar possibility was suggested by Dirac. He said: Let's admit that an electron acts on itself through the *second* term in Eq. (28.9) but not through the first. He then had an ingenious idea for getting rid of one but not the

other. Look, he said, we made a special assumption when we took only the *retarded* wave solutions of Maxwell's equations; if we were to take the *advanced* waves instead, we would get something different. The formula for the self-force would be

$$F = \alpha \frac{e^2}{ac^2} \ddot{x} + \frac{2}{3} \frac{e^2}{c^3} \ddot{x} + \gamma \frac{e^2 a}{c^4} \dot{x} \quad (28.12)$$

This equation is just like Eq. (28.9) except for the sign of the second term—and some higher terms—of the series [Changing from retarded to advanced waves is just changing the *sign* of the delay which, it is not hard to see, is equivalent to changing the sign of t everywhere. The only effect on Eq. (28.9) is to change the sign of all the odd time derivatives.] So, Dirac said, let's make the new rule that an electron acts on itself by one-half the *difference* of the retarded and advanced fields which it produces. The difference of Eqs. (28.9) and (28.12), divided by two, is then

$$F = -\frac{2}{3} \frac{e^2}{c^3} \ddot{x} + \text{higher terms.}$$

In all the higher terms, the radius a appears to some positive power in the numerator. Therefore, when we go to the limit of a point charge, we get only the one term—just what is needed. In this way, Dirac got the radiation resistance force and none of the inertial forces. There is no electromagnetic mass, and the classical theory is saved—but at the expense of an arbitrary assumption about the self-force.

The arbitrariness of the extra assumption of Dirac was removed, to some extent at least, by Wheeler and Feynman, who proposed a still stranger theory. They suggest that point charges interact *only* with other charges, but that the interaction is half through the advanced and half through the retarded waves. It turns out, most surprisingly, that in most situations you won't see any effects of the advanced waves, but they do have the effect of producing just the radiation reaction force. The radiation resistance is *not* due to the electron acting on itself, but from the following peculiar effect. When an electron is accelerated at the time t , it shakes all the other charges in the world at a *later* time $t' = t + r/c$ (where r is the distance to the other charge), because of the *retarded* waves. But then these other charges react back on the original electron through their *advanced* waves, which will arrive at the time t'' , equal to t' minus r/c , which is, of course, just t . (They also react back with their retarded waves too, but that just corresponds to the normal “reflected” waves.) The combination of the advanced and retarded waves means that at the instant it is accelerated an oscillating charge feels a force from all the charges that are “going to” absorb its radiated waves. You see what tight knots people have gotten into in trying to get a theory of the electron!

We'll describe now still another kind of theory, to show the kind of things that people think of when they are stuck. This is another modification of the laws of electrodynamics, proposed by Bopp. You realize that once you decide to change the equations of electromagnetism you can start anywhere you want. You can change the force law for an electron, or you can change the Maxwell equations (as we saw in the examples we have described), or you can make a change somewhere else. One possibility is to change the formulas that give the potentials in terms of the charges and currents. One of our formulas has been that the potentials at some point are given by the current density (or charge) at each other point at an earlier time. Using our four-vector notation for the potentials, we write

$$A_\mu(1, t) = \frac{1}{4\pi\epsilon_0 c^2} \int \frac{j_\mu(2, t - r_{12}/c)}{r_{12}} dV_2 \quad (28.13)$$

Bopp's beautifully simple idea is that: Maybe the trouble is in the $1/r$ factor in the integral. Suppose we were to start out by assuming only that the potential at one point depends on the charge density at any other point as *some* function of the distance between the points, say as $f(r_{12})$. The total potential at point (1)

will then be given by the integral of j_μ times this function over all space:

$$A_\mu(1) = \int j_\mu(2) f(r_{12}) dV_2.$$

That's all. No differential equation, nothing else. Well, one more thing. We also ask that the result should be relativistically invariant. So by "distance" we should take the invariant "distance" between two points in space-time. This distance squared (within a sign which doesn't matter) is

$$\begin{aligned} s_{12}^2 &= c^2(t_1 - t_2)^2 - r_{12}^2 \\ &= c^2(t_1 - t_2)^2 - (x_1 - x_2)^2 - (y_1 - y_2)^2 - (z_1 - z_2)^2. \end{aligned} \quad (28.14)$$

So, for a relativistically invariant theory, we should take some function of the magnitude of s_{12} , or what is the same thing, some function of s_{12}^2 . So Bopp's theory is that

$$A_\mu(1, t_1) = \int j_\mu(2, t_2) F(s_{12}^2) dV_2 dt_2. \quad (28.15)$$

(The integral must, of course, be over the four-dimensional volume $dt_2 dx_2 dy_2 dz_2$)

All that remains is to choose a suitable function for F . We assume only one thing about F —that it is very small except when its argument is near zero—so that a graph of F would be a curve like the one in Fig. 28-4. It is a narrow spike with a finite area centered at $s^2 = 0$, and with a width which we can say is roughly a^2 . We can say, crudely, that when we calculate the potential at point (1), only those points (2) produce any appreciable effect if $s_{12}^2 = c^2(t_2 - t_1)^2 - r_{12}^2$ is within $\pm a^2$ of zero. We can indicate this by saying that F is important only for

$$s_{12}^2 = c^2(t_1 - t_2)^2 - r_{12}^2 \approx \pm a^2. \quad (28.16)$$

You can make it more mathematical if you want to, but that's the idea.

Now suppose that a is very small in comparison with the size of ordinary objects like motors, generators, and the like so that for normal problems $r_{12} \gg a$. Then Eq. (28.16) says that charges contribute to the integral of Eq. (28.15) only when $t_1 - t_2$ is in the small range

$$c(t_1 - t_2) \approx \sqrt{r_{12}^2 \pm a^2} \approx r_{12} \sqrt{1 \pm \frac{a^2}{r_{12}^2}}.$$

Since $a^2/r_{12}^2 \ll 1$, the square root can be approximated by $1 \pm a^2/2r_{12}^2$, so

$$t_1 - t_2 = \frac{r_{12}}{c} \left(1 \pm \frac{a^2}{2r_{12}^2} \right) = \frac{r_{12}}{c} \pm \frac{a^2}{2r_{12}c}.$$

What is the significance? This result says that the only times t_2 that are important in the integral of A_μ are those which differ from the time t_1 , at which we want the potential, by the delay r_{12}/c —with a negligible correction so long as $r_{12} \gg a$. In other words, this theory of Bopp approaches the Maxwell theory—so long as we are far away from any particular charge—in the sense that it gives the retarded wave effects.

We can, in fact, see approximately what the integral of Eq. (28.15) is going to give. If we integrate first over t_2 from $-\infty$ to $+\infty$ —keeping r_{12} fixed—then s_{12}^2 is also going to go from $-\infty$ to $+\infty$. The integral will all come from t_2 's in a small interval of width $\Delta t_2 = 2 \times a^2/2r_{12}c$, centered at $t_1 - r_{12}/c$. Say that the function $F(s^2)$ has the value K at $s^2 = 0$; then the integral over t_2 gives approximately $K j_\mu \Delta t_2$, or

$$\frac{Ka^2}{c} \frac{j_\mu}{r_{12}}.$$

We should, of course, take the value of j_μ at $t_2 = t_1 - r_{12}/c$, so that Eq. (28.15) becomes

$$A_\mu(1, t_1) = \frac{Ka^2}{c} \int \frac{j_\mu(2, t_1 - r_{12}/c)}{r_{12}} dV_2.$$

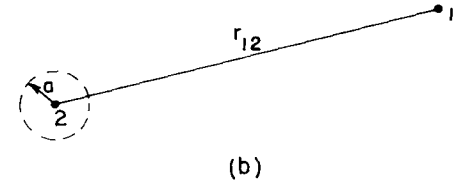
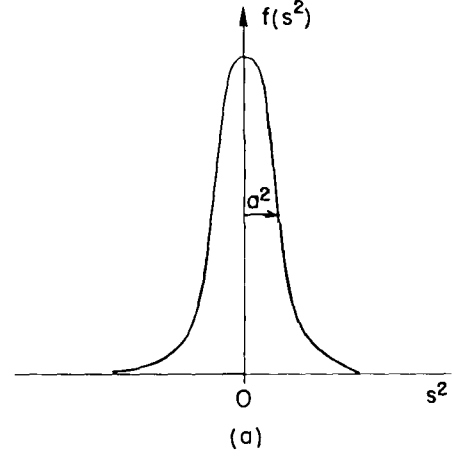


Fig. 28-4. The function $F(s^2)$ used in the nonlocal theory of Bopp.

If we pick $K = q^2 c / 4\pi\epsilon_0 a^2$, we are right back to the retarded potential solution of Maxwell's equations—including automatically the $1/r$ dependence! And it all came out of the simple proposition that the potential at one point in space-time depends on the current density at all other points in space-time, but with a weighting factor that is some narrow function of the four-dimensional distance between the two points. This theory again predicts a finite electromagnetic mass for the electron, and the energy and mass have the right relation for the relativity theory. They must, because the theory is relativistically invariant from the start, and everything seems to be all right.

There is, however, one fundamental objection to this theory and to all the other theories we have described. All particles we know obey the laws of quantum mechanics, so a quantum-mechanical modification of electrodynamics has to be made. Light behaves like photons. It isn't 100 percent like the Maxwell theory. So the electrodynamic theory has to be changed. We have already mentioned that it might be a waste of time to work so hard to straighten out the classical theory, because it could turn out that in quantum electrodynamics the difficulties will disappear or may be resolved in some other fashion. But the difficulties do not disappear in quantum electrodynamics. That is one of the reasons that people have spent so much effort trying to straighten out the classical difficulties, hoping that if they *could* straighten out the classical difficulty and *then* make the quantum modifications, everything would be straightened out. The Maxwell theory still has the difficulties after the quantum mechanics modifications are made.

The quantum effects do make some changes—the formula for the mass is modified, and Planck's constant \hbar appears—but the answer still comes out infinite unless you cut off an integration somehow—just as we had to stop the classical integrals at $r = a$. And the answers depend on how you stop the integrals. We cannot, unfortunately, demonstrate for you here that the difficulties are really basically the same, because we have developed so little of the theory of quantum mechanics and even less of quantum electrodynamics. So you must just take our word that the quantized theory of Maxwell's electrodynamics gives an infinite mass for a point electron.

It turns out, however, that nobody has ever succeeded in making a *self-consistent* quantum theory out of *any* of the modified theories. Born and Infeld's ideas have never been satisfactorily made into a quantum theory. The theories with the advanced and retarded waves of Dirac, or of Wheeler and Feynman, have never been made into a satisfactory quantum theory. The theory of Bopp has never been made into a satisfactory quantum theory. So today, there is no known solution to this problem. We do not know how to make a consistent theory—including the quantum mechanics—which does not produce an infinity for the self-energy of an electron, or any point charge. And at the same time, there is no satisfactory theory that describes a non-point charge. It's an unsolved problem.

In case you are deciding to rush off to make a theory in which the action of an electron on itself is completely removed, so that electromagnetic mass is no longer meaningful, and then to make a quantum theory of it, you should be warned that you are certain to be in trouble. There is definite experimental evidence of the existence of electromagnetic inertia—there is evidence that some of the mass of charged particles is electromagnetic in origin.

It used to be said in the older books that since Nature will obviously not present us with two particles—one neutral and the other charged, but otherwise the same—we will never be able to tell how much of the mass is electromagnetic and how much is mechanical. But it turns out that Nature *has* been kind enough to present us with just such objects, so that by comparing the observed mass of the charged one with the observed mass of the neutral one, we can tell whether there is any electromagnetic mass. For example, there are the neutrons and protons. They interact with tremendous forces—the nuclear forces—whose origin is unknown. However, as we have already described, the nuclear forces have one remarkable property. So far as they are concerned, the neutron and proton are exactly the same. The *nuclear* forces between neutron and neutron, neutron and proton, and proton and proton are all identical as far as we can tell. Only the little

electromagnetic forces are different; electrically the proton and neutron are as different as night and day. This is just what we wanted. There are two particles, identical from the point of view of the strong interactions, but different electrically. And they have a small difference in mass. The mass difference between the proton and the neutron—expressed as the difference in the rest-energy mc^2 in units of Mev—is about 1.3 Mev, which is about 2.6 times the electron mass. The classical theory would then predict a radius of about $\frac{1}{3}$ to $\frac{1}{2}$ the classical electron radius, or about 10^{-13} cm. Of course, one should really use the quantum theory, but by some strange accident, all the constants— 2π 's and \hbar 's, etc.—come out so that the quantum theory gives roughly the same radius as the classical theory. The only trouble is that the *sign* is wrong! The neutron is *heavier* than the proton.

Table 28-1
Particle Masses

Particle	Charge (electronic)	Mass (Mev)	Δm^* (Mev)
n (neutron)	0	939.5	
p (proton)	+1	938.2	-1.3
π (π -meson)	0	135.0	
	± 1	139.6	+4.6
K (K-meson)	0	497.8	
	± 1	493.9	-3.9
Σ (sigma)	0	1191.5	
	+1	1189.4	-2.1
	-1	1196.0	+4.5

* Δm = (mass of charged) - (mass of neutral).

Nature has also given us several other pairs—or triplets—of particles which appear to be exactly the same except for their electrical charge. They interact with protons and neutrons, through the so-called “strong” interactions of the nuclear forces. In such interactions, the particles of a given kind—say the π -mesons—behave in every way like one object *except* for their electrical charge. In Table 28-1 we give a list of such particles, together with their measured masses. The charged π -mesons—positive or negative—have a mass of 139.6 Mev, but the neutral π -meson is 4.6 Mev lighter. We believe that this mass difference is electromagnetic; it would correspond to a particle radius of 3 to 4×10^{-14} cm. You will see from the table that the mass differences of the other particles are usually of the same general size.

Now the size of these particles can be determined by other methods, for instance by the diameters they appear to have in high-energy collisions. So the electromagnetic mass seems to be in general agreement with electromagnetic theory, if we stop our integrals of the field energy at the same radius obtained by these other methods. That’s why we believe that the differences do represent electromagnetic mass.

You are no doubt worried about the different signs of the mass differences in the table. It is easy to see why the charged ones should be heavier than the neutral ones. But what about those pairs like the proton and the neutron, where the measured mass comes out the other way? Well, it turns out that these particles are complicated, and the computation of the electromagnetic mass must be more elaborate for them. For instance, although the neutron has no *net* charge, it *does* have a charge distribution inside it—it is only the *net* charge that is zero. In fact, we believe that the neutron looks—at least sometimes—like a proton with a negative π -meson in a “cloud” around it, as shown in Fig. 28-5. Although the neutron is “neutral,” because its total charge is zero, there are still electromagnetic energies

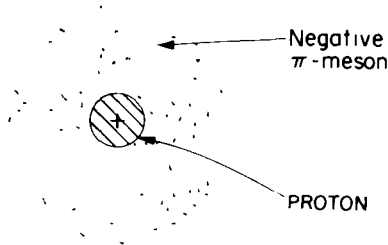


Fig. 28-5. A neutron may exist, at times, as a proton surrounded by a negative π -meson.

(for example, it has a magnetic moment), so it's not easy to tell the sign of the electromagnetic mass difference without a detailed theory of the internal structure.

We only wish to emphasize here the following points: (1) the electromagnetic theory predicts the existence of an electromagnetic mass, but it also falls on its face in doing so, because it does not produce a consistent theory—and the same is true with the quantum modifications; (2) there is experimental evidence for the existence of electromagnetic mass; and (3) all these masses are roughly the same as the mass of an electron. So we come back again to the original idea of Lorentz—maybe all the mass of an electron is purely electromagnetic, maybe the whole 0.511 Mev is due to electrodynamics. Is it or isn't it? We haven't got a theory, so we cannot say.

We must mention one more piece of information, which is the most annoying. There is another particle in the world called a *muon*—or μ -meson—which, so far as we can tell, differs in no way whatsoever from an electron except for its mass. It acts in every way like an electron: it interacts with neutrinos and with the electromagnetic field, and it has no nuclear forces. It does nothing different from what an electron does—at least, nothing which cannot be understood as merely a consequence of its higher mass (206.77 times the electron mass). Therefore, whenever someone finally gets the explanation of the mass of an electron, he will then have the puzzle of where a muon gets its mass. Why? Because whatever the electron does, the muon does the same—so the mass ought to come out the same. There are those who believe faithfully in the idea that the muon and the electron are the same particle and that, in the final theory of the mass, the formula for the mass will be a quadratic equation with two roots—one for each particle. There are also those who propose it will be a transcendental equation with an infinite number of roots, and who are engaged in guessing what the masses of the other particles in the series must be, and why these particles haven't been discovered yet.

28-6 The nuclear force field

We would like to make some further remarks about the part of the mass of nuclear particles that is not electromagnetic. Where does this other large fraction come from? There are other forces besides electrodynamics—like nuclear forces—that have their own field theories, although no one knows whether the current theories are right. These theories also predict a field energy which gives the nuclear particles a mass term analogous to electromagnetic mass; we could call it the “ π -mesic-field-mass.” It is presumably very large, because the forces are great, and it is the possible origin of the mass of the heavy particles. But the meson field theories are still in a most rudimentary state. Even with the well-developed theory of electromagnetism, we found it impossible to get beyond first base in explaining the electron mass. With the theory of the mesons, we strike out.

We may take a moment to outline the theory of the mesons, because of its interesting connection with electrodynamics. In electrodynamics, the field can be described in terms of a four-potential that satisfies the equation

$$\square^2 A_\mu = \text{sources.}$$

Now we have seen that pieces of the field can be radiated away so that they exist separated from the sources. These are the photons of light, and they are described by a differential equation without sources:

$$\square^2 A_\mu = 0.$$

People have argued that the field of nuclear forces ought also to have its own “photons”—they would presumably be the π -mesons—and that they should be described by an analogous differential equation. (Because of the weakness of the human brain, we can't think of something really new, so we argue by analogy with what we know.) So the meson equation might be

$$\square^2 \phi = 0,$$

where ϕ could be a different four-vector or perhaps a scalar. It turns out that the pion has no polarization, so ϕ should be a scalar. With the simple equation $\square^2 \phi = 0$, the meson field would vary with distance from a source as $1/r^2$, just as the electric field does. But we know that nuclear forces have much shorter distances of action, so the simple equation won't work. There is one way we can change things without disrupting the relativistic invariance: we can add or subtract from the D'Alembertian a constant, times ϕ . So Yukawa suggested that the free quanta of the nuclear force field might obey the equation

$$\square^2 \phi - \mu^2 \phi = 0, \quad (28.17)$$

where μ^2 is a constant—that is, an invariant scalar. (Since \square^2 is a scalar differential operator in four dimensions, its invariance is unchanged if we add another scalar to it.)

Let's see what Eq. (28.17) gives for the nuclear force when things are not changing with time. We want a spherically symmetric solution of

$$\nabla^2 \phi - \mu^2 \phi = 0$$

around some point source at, say, the origin. If ϕ depends only on r , we know that

$$\nabla^2 \phi = \frac{1}{r} \frac{\partial^2}{\partial r^2} (r\phi).$$

So we have the equation

$$\frac{1}{r} \frac{\partial^2}{\partial r^2} (r\phi) - \mu^2 \phi = 0$$

or

$$\frac{\partial^2}{\partial r^2} (r\phi) = \mu^2 (r\phi).$$

Thinking of $(r\phi)$ as our dependent variable, this is an equation we have seen many times. It's solution is

$$r\phi = Ke^{\pm\mu r}.$$

Clearly, ϕ cannot become infinite for large r , so the $+$ sign in the exponent is ruled out. The solution is

$$\phi = K \frac{e^{-\mu r}}{r}. \quad (28.18)$$

This function is called the *Yukawa potential*. For an attractive force, K is a negative number whose magnitude must be adjusted to fit the experimentally observed strength of the forces.

The Yukawa potential of the nuclear forces dies off more rapidly than $1/r$ by the exponential factor. The potential—and therefore the force—falls to zero much more rapidly than $1/r$ for distances beyond $1/\mu$, as shown in Fig. 28-6. The “range” of nuclear forces is much less than the “range” of electrostatic forces. It is found experimentally that the nuclear forces do not extend beyond about 10^{-13} cm, so $\mu \approx 10^{15} \text{ m}^{-1}$.

Finally, let's look at the free-wave solution of Eq. (28.17). If we substitute

$$\phi = \phi_0 e^{i(\omega t - kz)}$$

into Eq. (28.17), we get that

$$\frac{\omega^2}{c^2} - k^2 - \mu^2 = 0.$$

Relating frequency to energy and wave number to momentum, as we did at the end of Chapter 36 of Vol. I, we get that

$$\frac{E^2}{c^2} - p^2 = \mu^2 \hbar^2,$$

which says that the Yukawa “photon” has a mass equal to $\mu\hbar/c$. If we use for μ

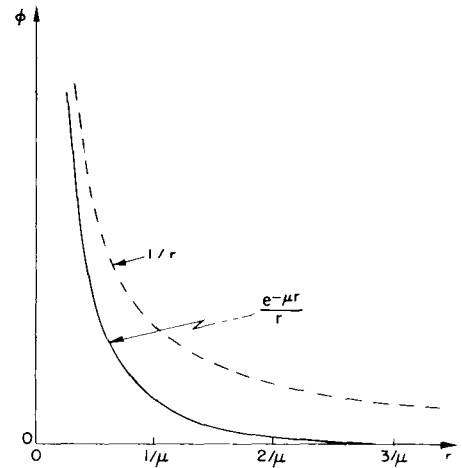


Fig. 28-6. The Yukawa potential $e^{-\mu r}/r$, compared with the Coulomb potential $1/r$.

the estimate 10^{15} m^{-1} , which gives the observed range of the nuclear forces, the mass comes out to $3 \times 10^{-25} \text{ gm}$, or 170 Mev, which is roughly the observed mass of the π -meson. So, by an analogy with electrodynamics, we would say that the π -meson is the “photon” of the nuclear force field. But now we have pushed the ideas of electrodynamics into regions where they may not really be valid—we have gone beyond electrodynamics to the problem of the nuclear forces.

The Motion of Charges in Electric and Magnetic Fields

29-1 Motion in a uniform electric or magnetic field

We want now to describe—mainly in a qualitative way—the motions of charges in various circumstances. Most of the interesting phenomena in which charges are moving in fields occur in very complicated situations, with many, many charges all interacting with each other. For instance, when an electromagnetic wave goes through a block of material or a plasma, billions and billions of charges are interacting with the wave and with each other. We will come to such problems later, but now we just want to discuss the much simpler problem of the motions of a single charge in a *given* field. We can then disregard all other charges—except, of course, those charges and currents which exist somewhere to produce the fields we will assume.

We should probably ask first about the motion of a particle in a uniform electric field. At low velocities, the motion is not particularly interesting—it is just a uniform acceleration in the direction of the field. However, if the particle picks up enough energy to become relativistic, then the motion gets more complicated. But we will leave the solution for that case for you to play with.

Next, we consider the motion in a uniform magnetic field with zero electric field. We have already solved this problem—one solution is that the particle goes in a circle. The magnetic force $q\mathbf{v} \times \mathbf{B}$ is always at right angles to the motion, so $d\mathbf{p}/dt$ is perpendicular to \mathbf{p} and has the magnitude $v p/R$, where R is the radius of the circle.

$$F = qvB = \frac{vp}{R}.$$

The radius of the circular orbit is then

$$R = \frac{p}{qB}. \quad (29.1)$$

That is only one possibility. If the particle has a component of its motion along the field direction, that motion is constant, since there can be no component of the magnetic force in the direction of the field. The general motion of a particle in a uniform magnetic field is a constant velocity parallel to \mathbf{B} and a circular motion at right angles to \mathbf{B} —the trajectory is a cylindrical helix (Fig. 29-1). The radius of the helix is given by Eq. (29.1) if we replace p by p_{\perp} , the component of momentum at right angles to the field.

29-2 Momentum analysis

A uniform magnetic field is often used in making a “momentum analyzer,” or “momentum spectrometer,” for high-energy charged particles. Suppose that charged particles are shot into a uniform magnetic field at the point A in Fig. 29-2(a), the magnetic field being perpendicular to the plane of the drawing. Each particle will go into an orbit which is a circle whose radius is proportional to its momentum. If all the particles enter perpendicular to the edge of the field, they will leave the field at a distance x (from A) which is proportional to their momentum p . A counter placed at some point such as C will detect only those particles whose momentum is in an interval Δp near the momentum $p = qBx/2$.

It is, of course, not necessary that the particles go through 180° before they are counted, but the so-called “ 180° spectrometer” has a special property. It is not

29-1 Motion in a uniform electric or magnetic field

29-2 Momentum analysis

29-3 An electrostatic lens

29-4 A magnetic lens

29-5 The electron microscope

29-6 Accelerator guide fields

29-7 Alternating-gradient focusing

29-8 Motion in crossed electric and magnetic fields

Review Chapter 30, Vol. I, *Diffraction*

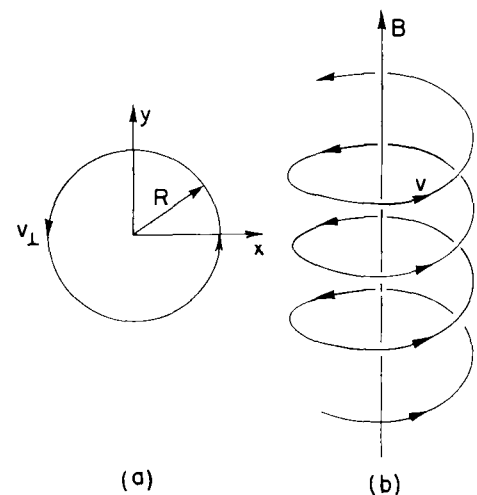


Fig. 29-1. Motion of a particle in a uniform magnetic field.

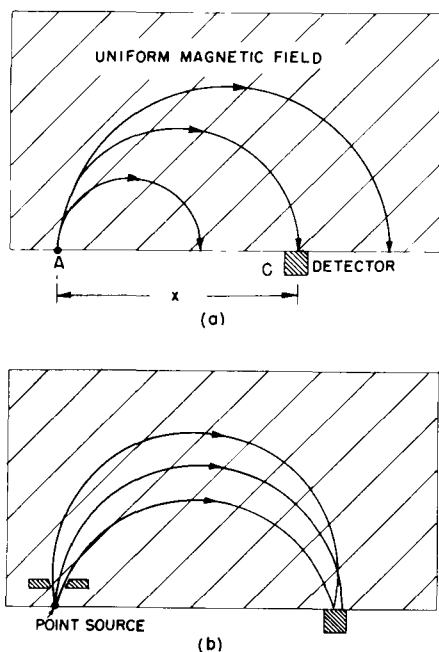


Fig. 29-2. A uniform-field, momentum spectrometer with 180° focusing: (a) different momenta; (b) different angles. (The magnetic field is directed perpendicular to the plane of the figure.)

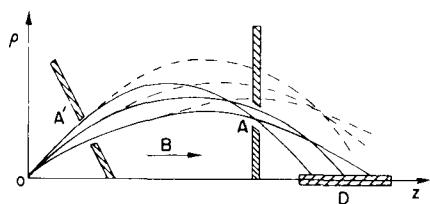


Fig. 29-3. An axial-field spectrometer.

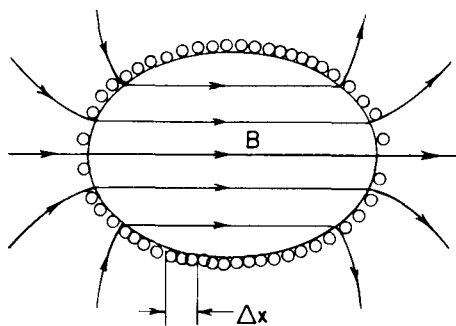


Fig. 29-4. An ellipsoidal coil with equal currents in each axial interval Δx produces a uniform magnetic field inside.

necessary that all the particles enter at right angles to the field edge. Figure 29-2(b) shows the trajectories of three particles, all with the *same* momentum but entering the field at different angles. You see that they take different trajectories, but all leave the field very close to the point C. We say that there is a “focus.” Such a focusing property has the advantage that larger angles can be accepted at A—although some limit is usually imposed, as shown in the figure. A larger angular acceptance usually means that more particles are counted in a given time, decreasing the time required for a given measurement.

By varying the magnetic field, or moving the counter along in x , or by using many counters to cover a range of x , the “spectrum” of momenta in the incoming beam can be measured. [By the “momentum spectrum” $f(p)$, we mean that the number of particles with momenta between p and $(p + dp)$ is $f(p) dp$.] Such measurements have been made, for example, to determine the distribution of energies in the β -decay of various nuclei.

There are many other forms of momentum spectrometers, but we will describe just one more, which has an especially large *solid* angle of acceptance. It is based on the helical orbits in a uniform field, like the one shown in Fig. 29-1. Let’s think of a cylindrical coordinate system— ρ , θ , z —set up with the z -axis along the direction of the field. If a particle is emitted from the origin at some angle α with respect to the z -axis, it will move along a spiral whose equation is

$$\rho = a \sin kz, \quad \theta = bz,$$

where a , b , and k are parameters you can easily work out in terms of p , α , and the magnetic field B . If we plot the distance ρ from the axis as a function of z for a given momentum, but for several starting angles, we will get curves like the solid ones drawn in Fig. 29-3. (Remember that this is just a kind of projection of a helical trajectory.) When the angle between the axis and the starting direction is larger, the peak value of ρ is large but the longitudinal velocity is less, so the trajectories for different angles tend to come to a kind of “focus” near the point A in the figure. If we put a narrow aperture of A , particles with a range of initial angles can still get through and pass on to the axis, where they can be counted by the long detector D .

Particles which leave the source at the origin with a higher momentum but at the same angles, follow the paths shown by the broken lines and do not get through the aperture at A . So the apparatus selects a small interval of momenta. The advantage over the first spectrometer described is that the aperture A —and the aperture A' —can be an annulus, so that particles which leave the source in a rather large solid angle are accepted. A large fraction of the particles from the source are used—an important advantage for weak sources or for very precise measurements.

One pays a price for this advantage, however, because a large volume of uniform magnetic field is required, and this is usually only practical for low-energy particles. One way of making a uniform field, you remember, is to wind a coil on a sphere, with a surface current density proportional to the sine of the angle. You can also show that the same thing is true for an ellipsoid of rotation. So such spectrometers are often made by winding an elliptical coil on a wooden (or aluminum) frame. All that is required is that the current in each interval of axial distance Δx be the same, as shown in Fig. 29-4

29-3 An electrostatic lens

Particle focusing has many applications. For instance, the electrons that leave the cathode in a TV picture tube are brought to a focus at the screen—to make a fine spot. In this case, one wants to take electrons all of the same energy but with different initial angles and bring them together in a small spot. The problem is like focusing light with a lens, and devices which do the corresponding job for particles are also called lenses.

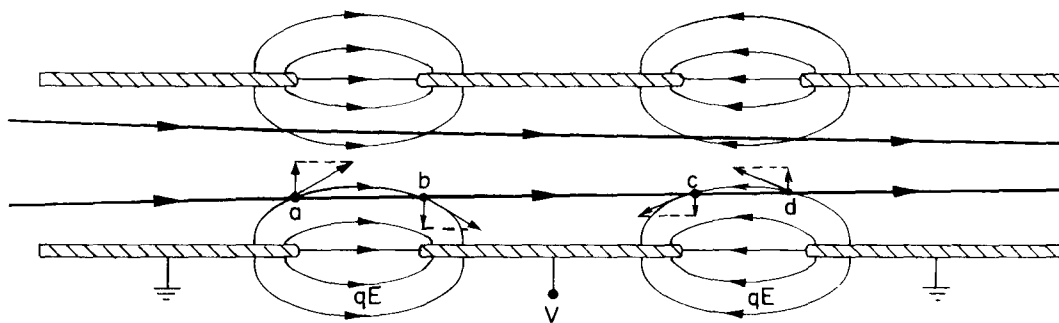


Fig. 29-5. An electrostatic lens. The field lines shown are "lines of force," that is, of qE .

One example of an electron lens is sketched in Fig. 29-5. It is an "electrostatic" lens whose operation depends on the electric field between two adjacent electrodes. Its operation can be understood by considering what happens to a parallel beam that enters from the left. When the electrons arrive at the region *a*, they feel a force with a sidewise component and get a certain impulse that bends them toward the axis. You might think that they would get an equal and opposite impulse in the region *b*, but that is not so. By the time the electrons reach *b* they have gained energy and so *spend less time* in the region *b*. The forces are the same, but the time is shorter, so the impulse is less. In going through the regions *a* and *b*, there is a net axial impulse, and the electrons are bent toward a common point. In leaving the high-voltage region, the particles get another kick toward the axis. The force is outward in region *c* and inward in region *d*, but the particles stay longer in the latter region, so there is again a net impulse. For distances not too far from the axis (Can you see why?), and this is just the condition necessary for lens-type focusing.

You can use the same arguments to show that there is focusing if the potential of the middle electrode is either positive or negative with respect to the other two. Electrostatic lenses of this type are commonly used in cathode-ray tubes and in some electron microscopes.

29-4 A magnetic lens

Another kind of lens—often found in electron microscopes—is the magnetic lens sketched schematically in Fig. 29-6. A cylindrically symmetric electromagnet has very sharp circular pole tips which produce a strong, nonuniform field in a small region. Electrons which travel vertically through this region are focused. You can understand the mechanism by looking at the magnified view of the pole-tip region drawn in Fig. 29-7. Consider two electrons *a* and *b* that leave the source *S* at some angle with respect to the axis. As electron *a* reaches the beginning of the field, it is deflected *away from you* by the horizontal component of the field. But then it will have a lateral velocity, so that when it passes through the strong vertical field, it will get an impulse toward the axis. Its lateral motion is taken out by the magnetic force as it leaves the field, so the net effect is an impulse toward the axis, plus a "rotation" about the axis. All the forces on particle *b* are opposite, so it also is deflected toward the axis. In the figure, the divergent electrons are brought into parallel paths. The action is like a lens with an object at the focal point. Another similar lens upstream can be used to focus the electrons back to a single point, making an image of the source *S*.

29-5 The electron microscope

You know that electron microscopes can "see" objects too small to be seen by optical microscopes. We discussed in Chapter 30 of Vol. I the basic limitations of any optical system due to diffraction of the lens opening. If a lens opening sub-

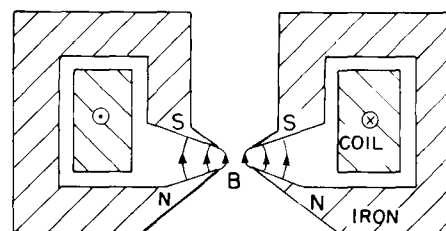


Fig. 29-6. A magnetic lens.

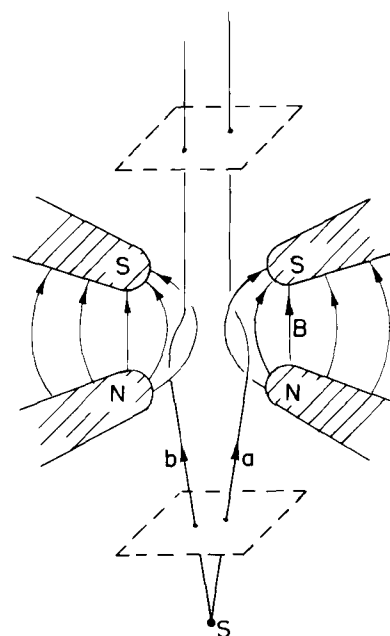


Fig. 29-7. Electron motion in the magnetic lens.

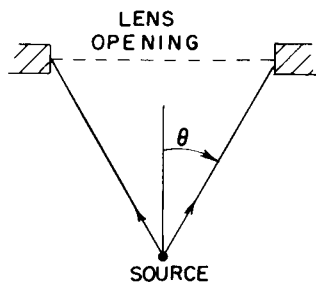


Fig. 29-8. The resolution of a microscope is limited by the angle subtended from the source.

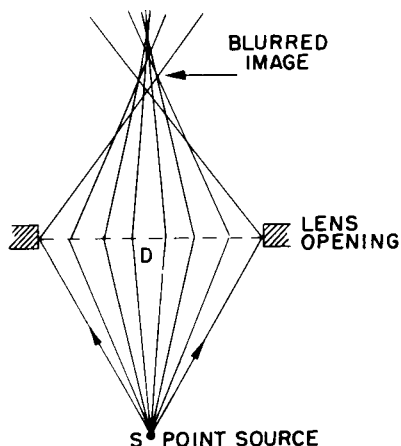


Fig. 29-9. Spherical aberration of a lens.

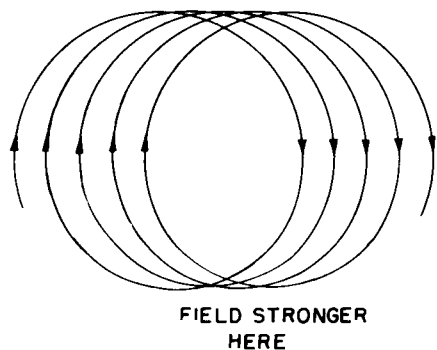


Fig. 29-10. Particle motion in a slightly nonuniform field.

tends the angle 2θ from a source (see Fig. 29-8), two neighboring spots at the source cannot be seen as separate if they are closer than about

$$\delta \approx \frac{\lambda}{\sin \theta},$$

where λ is the wavelength of the light. With the best optical microscope, θ approaches the theoretical limit of 90° , so δ is about equal to λ , or approximately 5000 angstroms.

The same limitation would also apply to an electron microscope, but there the wavelength is—for 50-kilovolt electrons—about 0.05 angstrom. If one could use a lens opening of near 30° , it would be possible to see objects only $\frac{1}{5}$ of an angstrom apart. Since the atoms in molecules are typically 1 or 2 angstroms apart, we could get photographs of molecules. Biology would be easy; we would have a photograph of the DNA structure. What a tremendous thing that would be! Most of present-day research in molecular biology is an attempt to figure out the shapes of complex organic molecules. If we could only see them!

Unfortunately, the best resolving power that has been achieved in an electron microscope is more like 20 angstroms. The reason is that no one has yet designed a lens with a large opening. All lenses have “spherical aberration,” which means that rays at large angles from the axis have a different point of focus than the rays nearer the axis, as shown in Fig. 29-9. By special techniques, optical microscope lenses can be made with a negligible spherical aberration, but no one has yet been able to make an electron lens which avoids spherical aberration.

In fact, one can show that any electrostatic or magnetic lens of the types we have described must have an irreducible amount of spherical aberration. This aberration—together with diffraction—limits the resolving power of electron microscopes to their present value.

The limitation we have mentioned does not apply to electric and magnetic fields which are not axially symmetric or which are not constant in time. Perhaps some day someone will think of a new kind of electron lens that will overcome the inherent aberration of the simple electron lens. Then we will be able to photograph atoms directly. Perhaps one day chemical compounds will be analyzed by looking at the positions of the atoms rather than by looking at the color of some precipitate!

29-6 Accelerator guide fields

Magnetic fields are also used to produce special particle trajectories in high-energy particle accelerators. Machines like the cyclotron and synchrotron bring particles to high energies by passing the particles repeatedly through a strong electric field. The particles are held in their cyclic orbits by a magnetic field.

We have seen that a particle in a uniform magnetic field will go in a circular orbit. This, however, is true only for a perfectly uniform field. Imagine a field B which is nearly uniform over a large area but which is slightly stronger in one region than in another. If we put a particle of momentum p in this field, it will go in a nearly circular orbit with the radius $R = p/qB$. The radius of curvature will, however, be slightly smaller in the region where the field is stronger. The orbit is not a closed circle but will “walk” through the field, as shown in Fig. 29-10. We can, if we wish, consider that the slight “error” in the field produces an extra angular kick which sends the particle off on a new track. If the particles are to make millions of revolutions in an accelerator, some kind of “radial focusing” is needed which will tend to keep the trajectories close to some design orbit.

Another difficulty with a uniform field is that the particles do not remain in a plane. If they start out with the slightest angle—or are given a slight angle by any small error in the field—they will go in a helical path that will eventually take them into the magnet pole or the ceiling or floor of the vacuum tank. Some arrangement must be made to inhibit such vertical drifts; the field must provide “vertical focusing” as well as radial focusing.

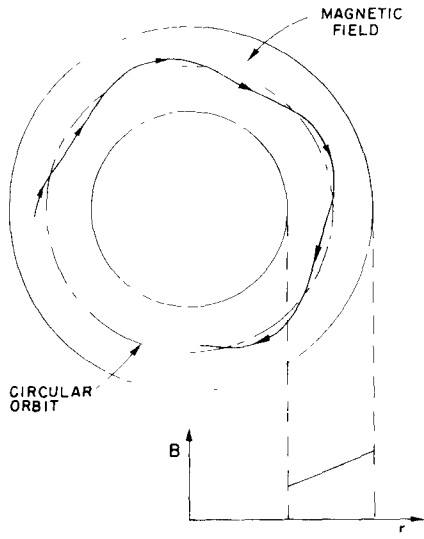


Fig. 29-11. Radial motion of a particle in a magnetic field with a large positive slope.

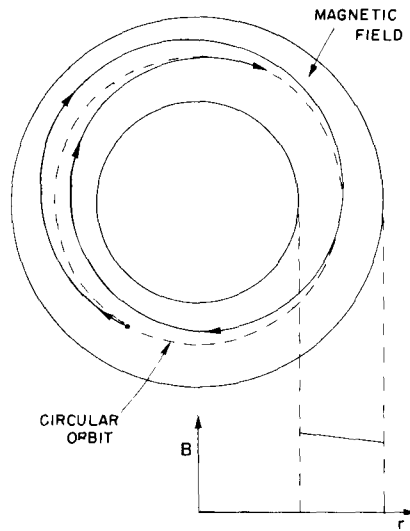


Fig. 29-12. Radial motion of a particle in a magnetic field with a small negative slope.

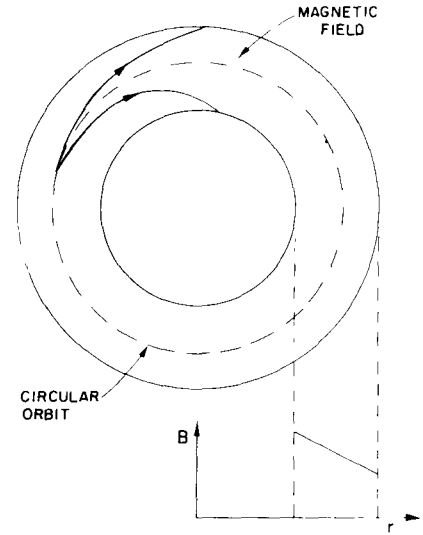


Fig. 29-13. Radial motion of a particle in a magnetic field with a large negative slope.

One would, at first, guess that radial focusing could be provided by making a magnetic field which increases with increasing distance from the center of the design path. Then if a particle goes out to a large radius, it will be in a stronger field which will bend it back toward the correct radius. If it goes to too small a radius, the bending will be less, and it will be returned toward the design radius. If a particle is once started at some angle with respect to the ideal circle, it will oscillate about the ideal circular orbit, as shown in Fig. 29-11. The radial focusing would keep the particles near the circular path.

Actually there is still some radial focusing even with the *opposite* field slope. This can happen if the radius of curvature of the trajectory does not increase more rapidly than the increase in the distance of the particle from the center of the field. The particle orbits will be as drawn in Fig. 29-12. If the gradient of the field is too large, however, the orbits will not return to the design radius but will spiral inward or outward, as shown in Fig. 29-13.

We usually describe the slope of the field in terms of the "relative gradient" or *field index*, n :

$$n = \frac{dB/B}{dr/r}. \quad (29.2)$$

A guide field gives radial focusing if this relative gradient is greater than -1 .

A radial field gradient will also produce *vertical* forces on the particles. Suppose we have a field that is stronger nearer to the center of the orbit and weaker at the outside. A vertical cross section of the magnet at right angles to the orbit might be as shown in Fig. 29-14. (For protons the orbits would be coming out of the page.) If the field is to be stronger to the left and weaker to the right, the lines of the magnetic field must be curved as shown. We can see that this must be so by using the law that the circulation of \mathbf{B} is zero in free space. If we take coordinates as shown in the figure, then

$$(\nabla \times \mathbf{B})_y = \frac{\partial B_x}{\partial z} - \frac{\partial B_z}{\partial x} = 0,$$

or

$$\frac{\partial B_x}{\partial z} = \frac{\partial B_z}{\partial x}. \quad (29.3)$$

Since we assume that $\partial B_z / \partial x$ is negative, there must be an equal negative $\partial B_x / \partial z$.

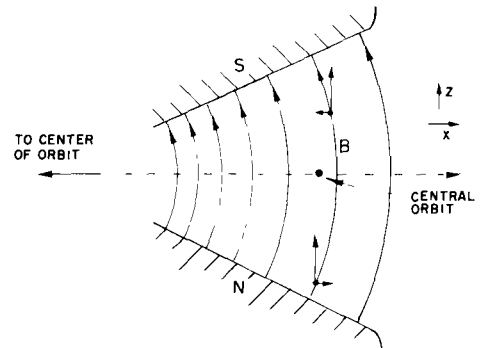


Fig. 29-14. A vertical guide field as seen in a cross section perpendicular to the orbits.

If the “nominal” plane of the orbit is a plane of symmetry where $B_r = 0$, then the radial component B_r will be negative above the plane and positive below. The lines must be curved as shown.

Such a field will have vertical focusing properties. Imagine a proton that is travelling more or less parallel to the central orbit but above it. The horizontal component of \mathbf{B} will exert a downward force on it. If the proton is below the central orbit, the force is reversed. So there is an effective “restoring force” toward the central orbit. From our arguments there will be vertical focusing, provided that the *vertical* field decreases with increasing radius; but if the field gradient is positive, there will be “vertical defocusing.” So for vertical focusing, the field index n must be less than zero. We found above that for radial focusing n had to be greater than -1 . The two conditions together give the condition that

$$-1 < n < 0$$

if the particles are to be kept in stable orbits. In cyclotrons, values very near zero are used; in betatrons and synchrotrons, the value $n = -0.6$ is typically used.

29-7 Alternating-gradient focusing

Such small values of n give rather “weak” focusing. It is clear that much more effective radial focusing would be given by a large positive gradient ($n \gg 1$), but then the vertical forces would be strongly defocusing. Similarly, large negative slopes ($n \ll -1$) would give stronger vertical forces but would cause radial defocusing. It was realized about 10 years ago, however, that a force that alternates between strong focusing and strong defocusing can still have a *net* focusing force.

To explain how *alternating-gradient focusing* works, we will first describe the operation of a quadrupole lens, which is based on the same principle. Imagine that a uniform negative magnetic field is added to the field of Fig. 29-14, with the strength adjusted to make zero field at the orbit. The resulting field—for small displacements from the neutral point—would be like the field shown in Fig. 29-15. Such a four-pole magnet is called a “quadrupole lens.” A positive particle that enters (from the reader) to the right or left of the center is pushed back toward the center. If the particle enters above or below, it is pushed *away* from the center. This is a horizontal focusing lens. If the horizontal gradient is reversed—as can be done by reversing all the polarities—the signs of all the forces are reversed and we have a vertical focusing lens, as in Fig. 29-16. For such lenses, the field strength—and therefore the focusing forces—increase linearly with the distance of the lens from the axis.

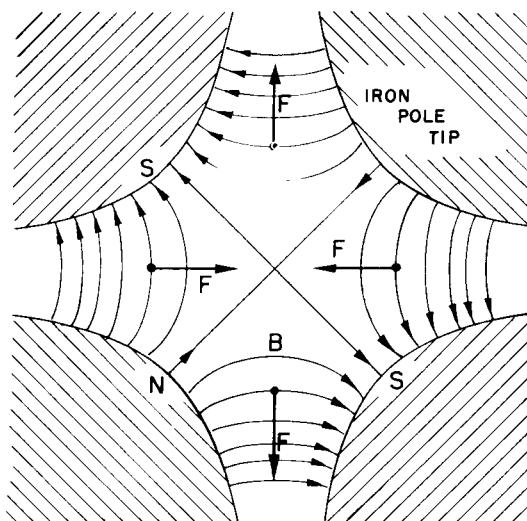


Fig. 29-15. A horizontal focusing quadrupole lens.

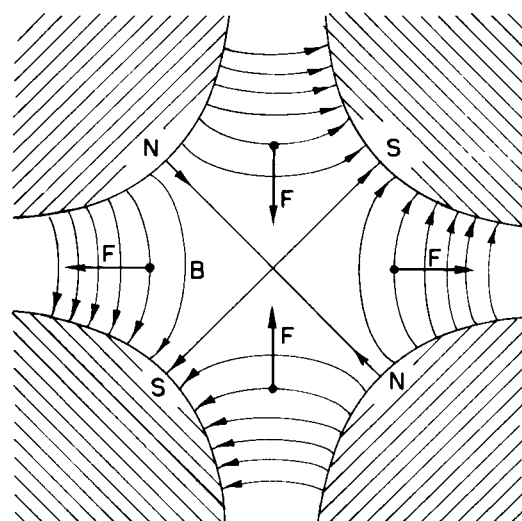


Fig. 29-16. A vertical focusing quadrupole lens.

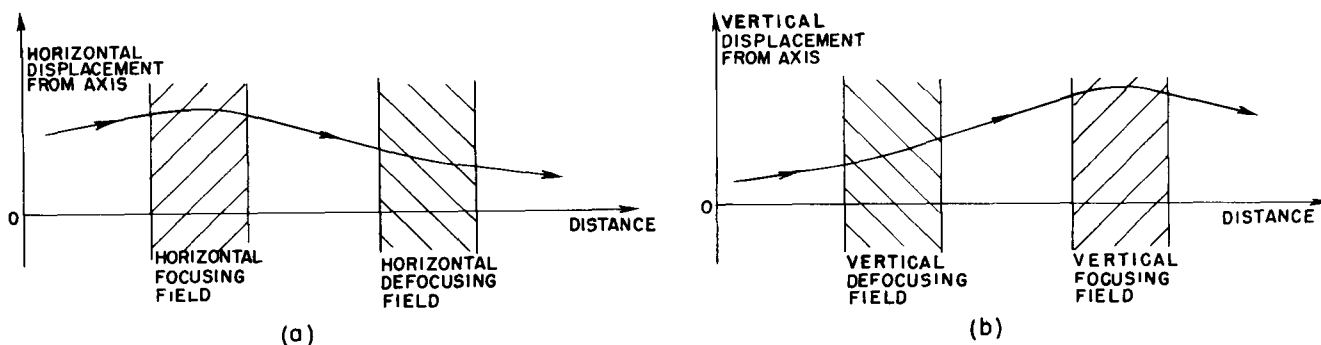


Fig. 29-17. Horizontal and vertical focusing with a pair of quadrupole lenses.

Now imagine that two such lenses are placed in series. If a particle enters with some horizontal displacement from the axis, as shown in Fig. 29-17(a), it will be deflected toward the axis in the first lens. When it arrives at the second lens it is closer to the axis, so the force outward is less and the outward deflection is less. There is a net bending toward the axis; the *average* effect is horizontally focusing. On the other hand, if we look at a particle which enters off the axis in the vertical direction, the path will be as shown in Fig. 29-17(b). The particle is first deflected *away* from the axis, but then it arrives at the second lens with a larger displacement, feels a stronger force, and so is bent toward the axis. Again the net effect is focusing. Thus a pair of quadrupole lenses acts independently for horizontal and vertical motion—very much like an optical lens. Quadrupole lenses are used to form and control beams of particles in much the same way that optical lenses are used for light beams.

We should point out that an alternating-gradient system does not *always* produce focusing. If the gradients are too large (in relation to the particle momentum or to the spacing between the lenses), the net effect can be a defocusing one. You can see how that could happen if you imagine that the spacing between the two lenses of Fig. 29-17 were increased, say, by a factor of three or four.

Let's return now to the synchrotron guide magnet. We can consider that it consists of an alternating sequence of "positive" and "negative" lenses with a superimposed uniform field. The uniform field serves to bend the particles, on the average, in a horizontal circle (with no effect on the vertical motion), and the alternating lenses act on any particles that might tend to go astray—pushing them always toward the central orbit (on the average).

There is a nice mechanical analog which demonstrates that a force which alternates between a "focusing" force and a "defocusing" force can have a net "focusing" effect. Imagine a mechanical "pendulum" which consists of a *solid* rod with a weight on the end, suspended from a pivot which is arranged to be moved rapidly up and down by a motor driven crank. Such a pendulum has *two* equilibrium positions. Besides the normal, downward-hanging position, the pendulum is also in equilibrium "hanging upward"—with its "bob" *above* the pivot! Such a pendulum is drawn in Fig. 29-18.

By the following argument you can see that the vertical pivot motion is equivalent to an alternating focusing force. When the pivot is accelerated downward, the "bob" tends to move inward, as indicated in Fig. 29-19. When the pivot is accelerated upward, the effect is reversed. The force restoring the "bob" toward the axis alternates, but the average effect is a force toward the axis. So the pendulum will swing back and forth about a neutral position which is just opposite the normal one.

There is, of course, a much easier way of keeping a pendulum upside down, and that is by *balancing* it on your finger! But try to balance *two independent* sticks on the *same finger*! Or one stick with your eyes closed! Balancing involves making a correction for what is going wrong. And this is not possible, in general, if there are several things going wrong at once. In a synchrotron there are billions of particles going around together, each one of which may start out with a different "error." The kind of focusing we have been describing works on them all.

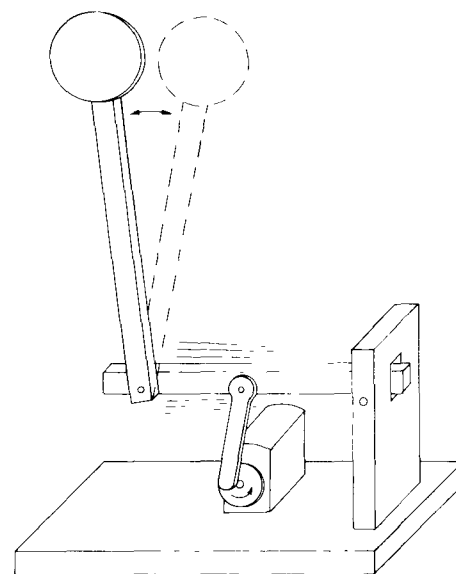


Fig. 29-18. A pendulum with an oscillating pivot can have a stable position with the bob above the pivot.

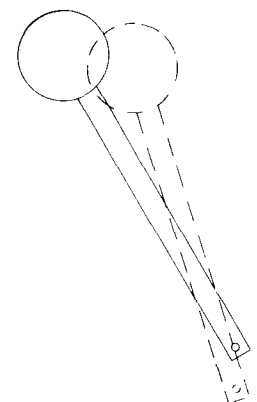


Fig. 29-19. A downward acceleration of the pivot causes the pendulum to move toward the vertical.

29-8 Motion in crossed electric and magnetic fields

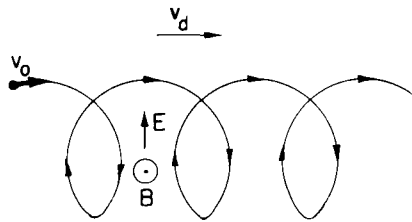


Fig. 29-20. Path of a particle in crossed electric and magnetic fields.

So far we have talked about particles in electric fields only or in magnetic fields only. There are some interesting effects when there are both kinds of fields at the same time. Suppose we have a uniform magnetic field \mathbf{B} and an electric field \mathbf{E} at right angles. Particles that start out perpendicular to \mathbf{B} will move in a curve like the one in Fig 29-20 (The figure is a *plane* curve, *not* a helix!) We can understand this motion qualitatively. When the particle (assumed positive) moves in the direction of \mathbf{E} , it picks up speed, and so it is bent less by the magnetic field. When it is going against the \mathbf{E} -field, it loses speed and is continually bent more by the magnetic field. The net effect is that it has an average “drift” in the direction of $\mathbf{E} \times \mathbf{B}$.

We can, in fact, show that the motion is a uniform circular motion superimposed on a uniform sidewise motion at the speed $v_d = E/B$ —the trajectory in Fig. 29-20 is a cycloid. Imagine an observer who is moving to the right at a constant speed. In his frame our magnetic field gets transformed to a new magnetic field *plus* an electric field in the *downward* direction. If he has just the right speed, his total electric field will be zero, and he will see the electron going in a circle. So the motion *we* see is a circular motion, plus a translation at the drift speed $v_d = E/B$. The motion of electrons in crossed electric and magnetic fields is the basis of the *magnetron* tubes, i.e., oscillators used for generating microwave energy.

There are many other interesting examples of particle motions in electric and magnetic fields—such as the orbits of the electrons and protons trapped in the Van Allen belts—but we do not, unfortunately, have the time to deal with them here

The Internal Geometry of Crystals

30-1 The internal geometry of crystals

We have finished the study of the basic laws of electricity and magnetism, and we are now going to study the electromagnetic properties of matter. We begin by describing solids—that is, crystals. When the atoms of matter are not moving around very much, they get stuck together and arrange themselves in a configuration with as low an energy as possible. If the atoms in a certain place have found a pattern which seems to be of low energy, then the atoms somewhere else will probably make the same arrangement. For these reasons, we have in a solid material a repetitive pattern of atoms

In other words, the conditions in a crystal are this way: The environment of a particular atom in a crystal has a certain arrangement, and if you look at the same kind of an atom at another place farther along, you will find one whose surroundings are exactly the same. If you pick an atom farther along by the same distance, you will find the conditions exactly the same once more. The pattern is repeated over and over again—and, of course, in three dimensions.

Imagine the problem of designing a wallpaper—or a cloth, or some geometric design for a plane area—in which you are supposed to have a design element which repeats and repeats and repeats, so that you can make the area as large as you want. This is the two-dimensional analog of a problem which a crystal solves in three dimensions. For example, Fig. 30-1(a) shows a common kind of wallpaper design. There is a single element repeated in a pattern that can go on forever. The geometric characteristics of this wallpaper design, considering only its repetition properties and not worrying about the geometry of the flower itself or its artistic merit, are contained in Fig. 30-1(b). If you start at any point, you can find the *corresponding* point by moving the distance a along the direction of arrow 1. You can also get to a corresponding point if you move the distance b in the direction of the other arrow. There are, of course, many other directions. You can go, for example, from point α to point β and reach a corresponding position, but such a step can be considered as a combination of a step along direction 1, followed by a step along direction 2. One of the basic properties of the pattern can be described by the two shortest steps to nearby equal positions. By “equal” positions we mean that if you were to stand in any one of them and look around you, you would see exactly the same thing as if you were to stand in another one. That’s the fundamental property of a crystal. The only difference is that a crystal is a three-dimensional arrangement instead of a two-dimensional arrangement; and naturally, instead of flowers, each element of the lattice is some kind of an arrangement of atoms—perhaps six hydrogen atoms and two carbon atoms—in some kind of pattern. The pattern of atoms in a crystal can be found out experimentally by x-ray diffraction. We have mentioned this method briefly before, and won’t say any more now except that the precise arrangement of the atoms in space has been worked out for most simple crystals and also for some fairly complex ones.

The internal pattern of a crystal shows up in several ways. First, the binding strength of the atoms in certain directions is usually stronger than in other directions. This means that there are certain planes through the crystal where it is more easily broken than others. They are called the *cleavage planes*. If you crack a crystal with a knife blade it will often split apart along such a plane. Second, the internal structure often appears at the surface because of the way the crystal was formed. Imagine a crystal being deposited out of a solution. There are the atoms floating around in the solution and finally settling down when they find a position

30-1 The internal geometry of crystals

30-2 Chemical bonds in crystals

30-3 The growth of crystals

30-4 Crystal lattices

30-5 Symmetries in two dimensions

30-6 Symmetries in three dimensions

30-7 The strength of metals

30-8 Dislocations and crystal growth

30-9 The Bragg-Nye crystal model

Reference: C. Kittel, *Introduction to Solid State Physics*, John Wiley and Sons, Inc., New York, 2nd ed., 1956.

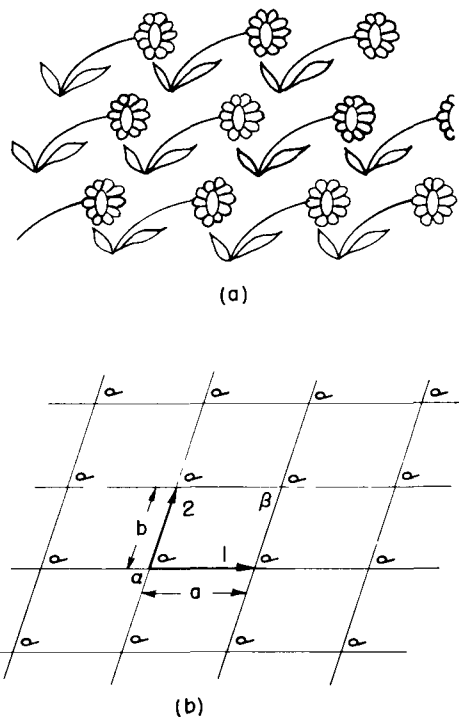


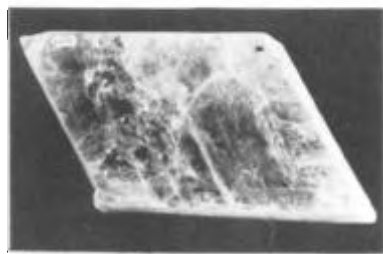
Fig. 30-1. A repeating pattern in two dimensions.



(a)



(b)



(c)

Fig. 30-2. Natural crystals: (a) quartz, (b) sodium chloride, (c) mica.

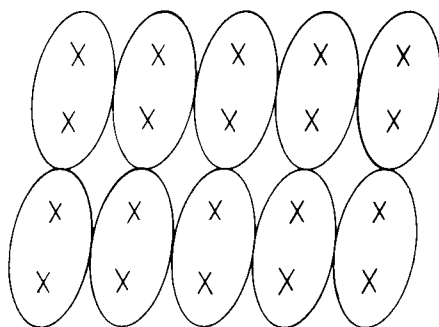


Fig. 30-3. The lattice of a molecular crystal.

of lowest energy. (It's as if the wallpaper got made by flowers drifting around until one drifted accidentally into place and got stuck, and then the next, and the next so that the pattern gradually grows.) You can appreciate that there will be certain directions in which it will grow at a different speed than in other directions, thereby growing into some kind of geometrical shape. Because of such effects, the outside surfaces of many crystals show some of the character of the internal arrangement of the atoms

For example, Fig. 30-2(a) shows the shape of a typical quartz crystal whose internal pattern is hexagonal. If you look closely at such a crystal, you will notice that the outside does not make a very good hexagon because the sides are not all of equal length—they are, in fact, often very unequal. But in one respect it is a very good hexagon: the *angles* between the faces are exactly 120° . Clearly, the size of any particular face is an accident of the growth, but the *angles* are a representation of the internal geometry. So every crystal of quartz has a different shape, even though the angles between corresponding faces are always the same.

The internal geometry of a crystal of sodium chloride is also evident from its external shape. Figure 30-2(b) shows the shape of a typical grain of salt. Again the crystal is not a perfect cube, but the faces *are* exactly at right angles to one another.

A more complicated crystal is mica, which has the shape shown in Fig. 30-2(c). It is a highly anisotropic crystal, as is easily seen from the fact that it is very tough if you try to pull it apart in one direction (horizontally in the figure), but very easy to split by pulling apart in the other direction (vertically). It has commonly been used to obtain very tough, thin sheets. Mica and quartz are two examples of natural minerals containing silica. A third example of a mineral with silica is asbestos, which has the interesting property that it is easily pulled apart in two directions but not in the third. It appears to be made of very strong, *linear* fibers.

30-2 Chemical bonds in crystals

The mechanical properties of crystals clearly depend on the kind of chemical bindings between the atoms. The strikingly different strength of mica along different directions depends on the kinds of interatomic binding in the different directions. You have already learned in chemistry, no doubt, about the different kinds of chemical bonds. First, there are ionic bonds, as we have already discussed for sodium chloride. Roughly speaking, the sodium atoms have lost an electron and become positive ions, the chlorine atoms have gained an electron and become negative ions. The positive and negative ions are arranged in a three-dimensional checkerboard and are held together by electrical forces.

The covalent bond—in which electrons are shared between two atoms—is more common and is usually very strong. In a diamond, for example, the carbon atoms have covalent bonds in all four directions to the nearest neighbors, so the crystal is very hard indeed. There is also covalent bonding between silicon and oxygen in a quartz crystal, but there the bond is really only partially covalent. Because there is not complete sharing of the electrons, the atoms are partly charged, and the crystal is somewhat ionic. Nature is not as simple as we try to make it; there are really all possible gradations between covalent and ionic bonding.

A sugar crystal has still another kind of binding. In it there are large molecules in which the atoms are held strongly together by covalent bonds, so that the molecule is a tough structure. But since the strong bonds are completely satisfied, there are only relatively weak attractions between the separate, individual molecules. In such *molecular* crystals the molecules keep their individual identity, so to speak, and the internal arrangement might be as shown in Fig. 30-3. Since the molecules are not held strongly to each other, the crystals are easy to break. They are quite different from something like diamond, which is really one giant molecule that cannot be broken anywhere without disrupting strong covalent bonds. Paraffin is another example of a molecular crystal.

An extreme example of a molecular crystal occurs in a substance like solid argon. There is very little attraction between the atoms—each atom is a completely

saturated monatomic molecule. But at very low temperatures, the thermal motion is very small, so the slight interatomic forces can cause the atoms to settle down into a regular array like a pile of closely packed spheres.

The metals form a completely different class of substances. The bonding is of an entirely different kind. In a metal the bonding is not between adjacent atoms but is a property of the whole crystal. The valence electrons are not attached to one atom or to a pair of atoms but are shared throughout the crystal. Each atom contributes an electron to a universal pool of electrons, and the atomic positive ions reside in the sea of negative electrons. The electron sea holds the ions together like some kind of glue.

In the metals, since there are no special bonds in any particular direction, there is no strong directionality in the binding. They are still crystalline, however, because the total energy is lowest when the atomic ions are arranged in some definite array—although the energy of the preferred arrangement is not usually much lower than other possible ones. To a first approximation, the atoms of many metals are like small spheres packed in as tightly as possible.

30-3 The growth of crystals

Try to imagine the natural formation of crystals in the earth. In the earth's surface there is a big mixture of all kinds of atoms. They are being continually churned about by volcanic action, by wind, and by water—continually being moved about and mixed. Yet, by some trick, silicon atoms gradually begin to find each other, and to find oxygen atoms, to make silica. One atom at a time is added to the others to build up a crystal—the mixture gets unmixed. And somewhere nearby, sodium and chlorine atoms are finding each other and building up a crystal of salt.

How does it happen that once a crystal is started, it permits only a particular kind of atom to join on? It happens because the whole system is working toward the lowest possible energy. A growing crystal will accept a new atom if it is going to make the energy as low as possible. But how does it *know* that a silicon—or an oxygen—atom at some particular spot is going to result in the lowest possible energy? It does it by trial and error. In the liquid, all of the atoms are in perpetual motion. Each atom bounces against its neighbors about 10^{13} times every second. If it hits against the right spot of growing crystal, it has a somewhat smaller chance of jumping off again if the energy is low. By continually testing over periods of millions of years at a rate of 10^{13} tests per second, the atoms gradually build up at the places where they find their lowest energy. Eventually they grow into big crystals.

30-4 Crystal lattices

The arrangement of the atoms in a crystal—the crystal *lattice*—can take on many geometric forms. We would like to describe first the simplest lattices, which are characteristic of most of the metals and of the solid form of the inert gases. They are the cubic lattices which can occur in two forms: the body-centered cubic, shown in Fig. 30-4(a), and the face-centered cubic shown in Fig. 30-4(b). The drawings show, of course, only one cube of the lattice; you are to imagine that the pattern is repeated indefinitely in three dimensions. Also, to make the drawing clearer, only the “centers” of the atoms are shown. In an actual crystal, the atoms are more like spheres in contact with each other. The dark and light spheres in the drawings may, in general, stand for different kinds of atoms or may be the same kind. For instance, iron has a body-centered cubic lattice at low temperatures, but a face-centered cubic lattice at higher temperatures. The physical properties are quite different in the two crystalline forms.

How do such forms come about? Imagine that you have the problem of packing spherical atoms together as tightly as possible. One way would be to start by making a layer in a “hexagonal close-packed array,” as shown in Fig. 30-5(a). Then you could build up a second layer like the first, but displaced horizontally,

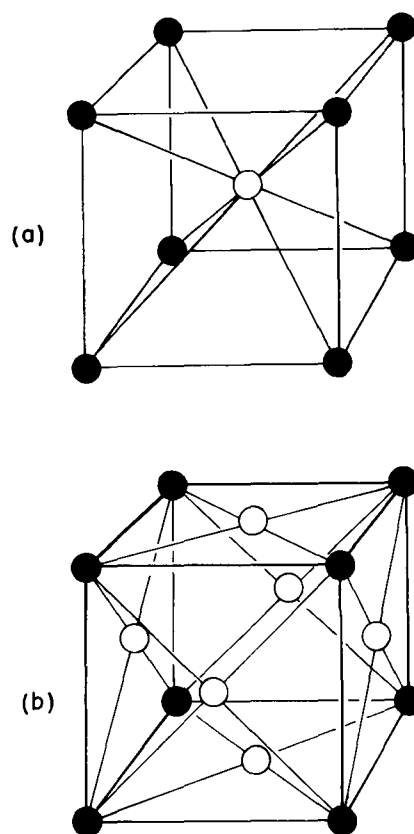


Fig. 30-4. The unit cell of cubic crystals: (a) body-centered, (b) face-centered.

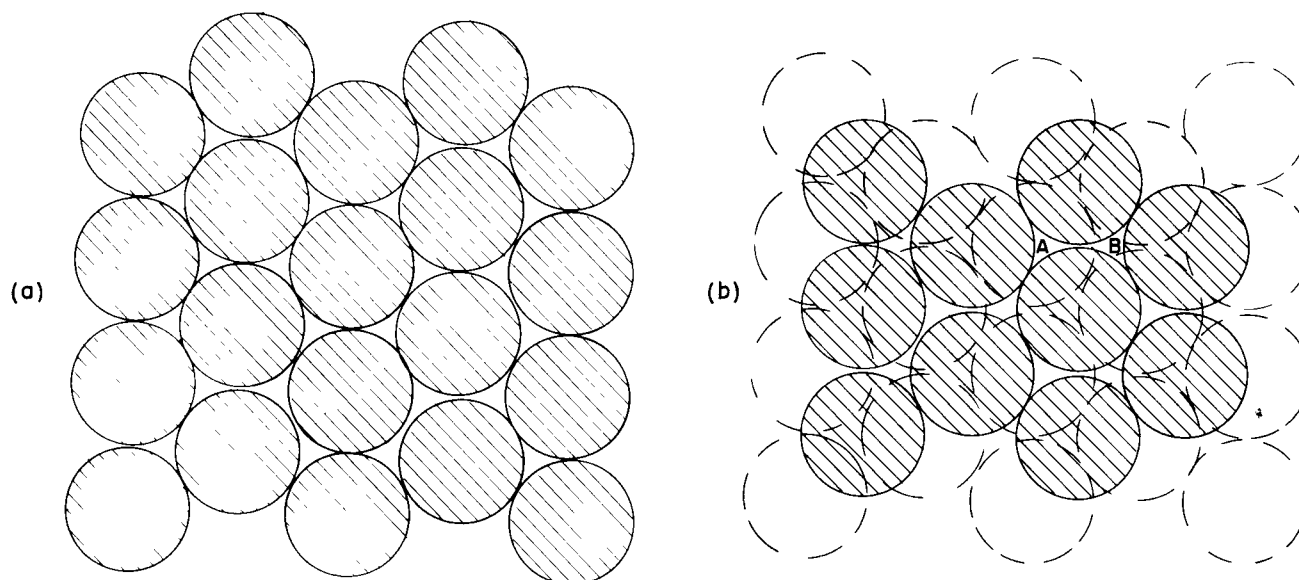


Fig. 30-5. Building up a hexagonal close-packed lattice.

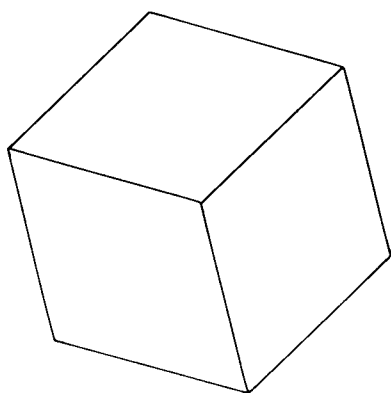


Fig. 30-6. Is this a hexagon or a cube seen from one corner?

as shown in Fig. 30-5(b). Next, you can put on the third layer. But notice! There are *two* distinct ways of placing the *third* layer. If you start the third layer by placing an atom at *A* in Fig. 30-5(b), each atom in the third layer is directly above an atom of the bottom layer. On the other hand, if you start the third layer by putting an atom at the position *B*, the atoms of the third layer will be centered at points exactly in the middle of a triangle formed by three atoms of the bottom layer. Any other starting place is equivalent to *A* or *B*, so there are only two ways of placing the third layer.

If the third layer has an atom at point *B*, the crystal lattice is a face-centered cubic—but seen at an angle. It seems funny that starting with hexagons you can end up with cubes. But notice that a cube looked at from a corner has a hexagonal outline. For instance, Fig. 30-6 could represent a plane hexagon or a cube seen in perspective!

If a third layer is added to Fig. 30-5(b) by starting with an atom at *A*, there is no cubical structure, and the lattice has instead only a hexagonal symmetry. It is clear that both possibilities we have described are equally close-packed.

Some metals—for example, copper and silver—choose the first alternative, the face-centered cubic. Others—for example, beryllium and magnesium—choose the other alternatives; they form hexagonal crystals. Clearly, which crystal lattice appears cannot depend only on the packing of little spheres, but must also be determined in part by other factors. In particular, it depends on the slight remaining angular dependence of the interatomic forces (or, in the case of the metals, on the energy of the electron pool). You will, no doubt, learn all about such things in your chemistry courses.

30-5 Symmetries in two dimensions

We would now like to discuss some of the properties of crystals from the point of view of their internal symmetries. The main feature of a crystal is that if you start at one atom and move to a corresponding atom one lattice unit away, you are again in the same kind of an environment. That's the fundamental proposition. But if you were an atom, there would be another kind of change that could take you again to the same environment—that is, another possible "symmetry." Figure 30-7(a) shows another possible "wallpaper-type" design (though one you have probably never seen). Suppose we compare the environments for points *A* and *B*. You might, at first, think that they are the same—but not quite. Points *C* and *D* are equivalent to *A*, but the environment of *B* is like that of *A* only if the surroundings are reversed, as in a mirror reflection.

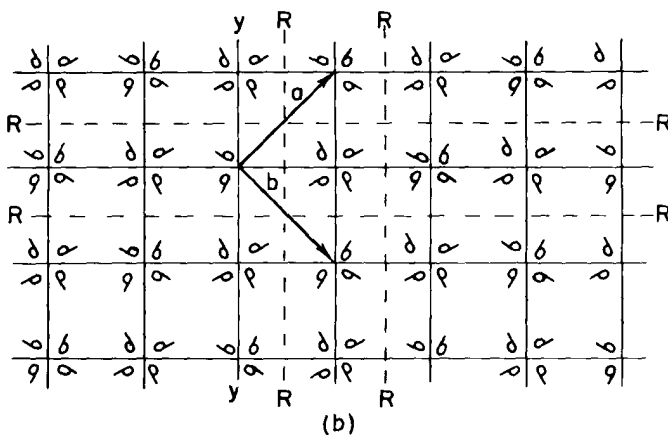
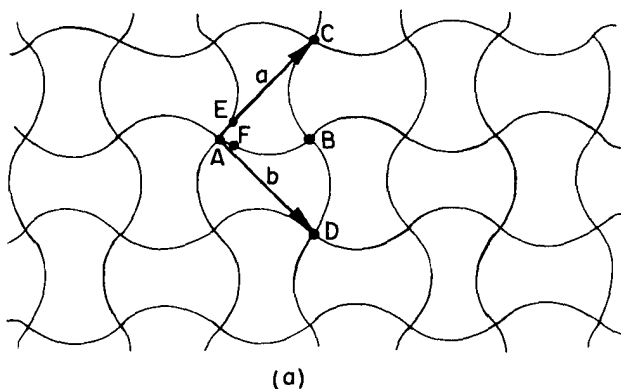


Fig. 30-7. A pattern of high symmetry.

There are other kinds of “equivalent” points in the pattern. For instance, the points *E* and *F* have the “same” environments except that one is rotated 90° with respect to the other. The pattern is quite special. A rotation of 90° —or any multiple of it—about a vertex such as *A* gives the same pattern all over again. A crystal with such a structure would have square corners on the outside, but inside it is more complicated than a simple cube.

Now that we have described some special examples, let’s try to figure out all the possible symmetries a crystal can have. First, we consider what happens in a plane. A *plane* lattice can be defined by the two so-called *primitive* vectors that go from one point of the lattice to the two *nearest* equivalent points. The two vectors **1** and **2** are the primitive vectors of the lattice of Fig. 30-1. The two vectors *a* and *b* of Fig. 30-7(a) are the primitive vectors of the pattern there. We could, of course, equally well replace *a* by $-a$, or *b* by $-b$. Since *a* and *b* are equal in magnitude and at right angles, a rotation of 90° turns *a* into *b*, and *b* into $-a$, giving the same lattice once again.

We see that there are lattices which have a “four-sided” symmetry. And we have described earlier a close-packed array based on a hexagon which could have a six-sided symmetry. A rotation of the array of circles in Fig. 30-5(a) by an angle of 60° about the center of any circle brings the pattern back to itself.

What other kinds of rotational symmetry are there? Can we have, for example, a fivefold or an eightfold rotational symmetry? It is easy to see that they are impossible. *The only symmetry with more sides than four is a six-sided symmetry.* First, let’s show that more than sixfold symmetry is impossible. Suppose we try to imagine a lattice with two equal primitive vectors with an enclosed angle less than 60° , as in Fig. 30-8(a). We are to suppose that points *B* and *C* are equivalent to *A*, and that *a* and *b* are the two *shortest* vectors from *A* to its equivalent neighbors. But that is clearly wrong, because the distance between *B* and *C* is shorter than from either one to *A*. There must be a neighbor at *D* equivalent to *A* which is closer than *B* or *C*. We should have chosen *b'* as one of our primitive vectors. So the angle between the two primitive vectors must be 60° or larger. Octagonal symmetry is not possible.

What about fivefold symmetry? If we assume that the primitive vectors *a* and *b* have equal lengths and make an angle of $2\pi/5 = 72^\circ$, as in Fig. 30-8(b), then there should also be an equivalent lattice point at *D*, at 72° from *C*. But the vector *b'* from *E* to *D* is then less than *b*, so *b* is not a primitive vector. There can be no fivefold symmetry. The only possibilities that do not get us into this kind of difficulty are $\theta = 60^\circ, 90^\circ$, or 120° . Zero or 180° are also clearly possible. One way of stating our result is that the pattern can be left unchanged by a rotation of one full turn (no change at all), one-half of a turn, one-third, one-fourth, or one-sixth of a turn. And those are all the possible rotational symmetries in a plane—a total of five. If $\theta = 2\pi/n$, we speak of an “*n*-fold” symmetry. We say

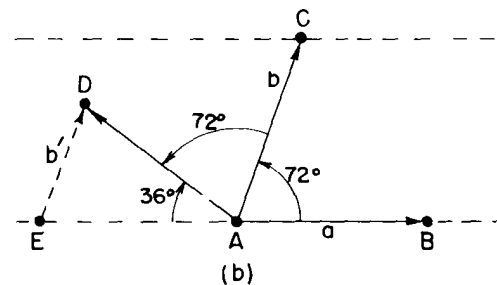
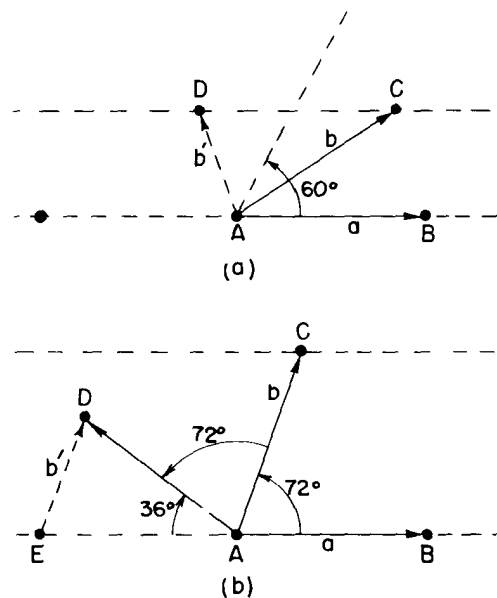


Fig. 30-8. (a) Rotational symmetries greater than sixfold are not possible. (b) Fivefold rotational symmetry is not possible.

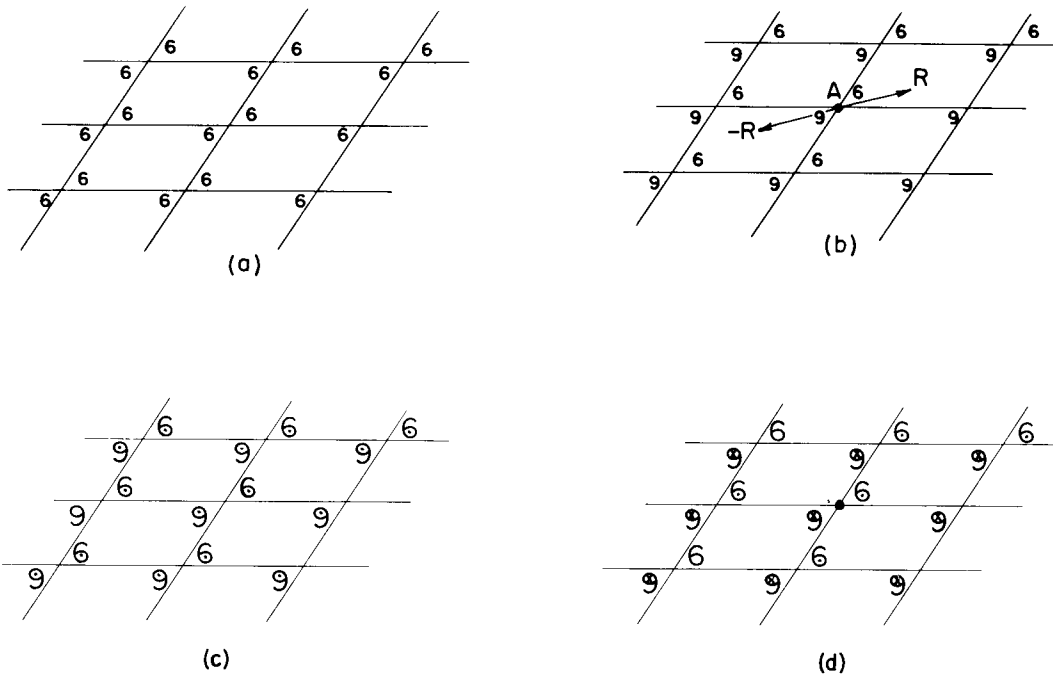


Fig. 30-9. Symmetry under inversion. Pattern (b) is unchanged if $R \rightarrow -R$, but pattern (a) is changed. In three dimensions pattern (d) is symmetric under an inversion but (c) is not.

that a pattern with n equal to 4 or to 6 has a “higher symmetry” than one with n equal to 1 or to 2.

Returning to Fig. 30-7(a), we see that the pattern has a fourfold rotational symmetry. We have drawn in Fig. 30-7(b) another design which has the same symmetry properties as part (a). The little comma-like figures are asymmetric objects which serve to define the symmetry of the design inside of each square. Notice that the commas are reversed in alternate squares, so that the unit cell is larger than one of the small squares. If there were no commas, the pattern would still have fourfold symmetry, but the unit cell would be smaller. The patterns of Fig. 30-7 also have other symmetry properties. For instance, a reflection about any of the broken lines $R-R$ reproduces the same pattern.

The patterns of Fig. 30-7 have still another kind of symmetry. If the pattern is reflected about the line $Y-Y$ and shifted one square to the right (or left), we get back the original pattern. The line $Y-Y$ is called a “glide” line.

These are all the possible symmetries in two dimensions. There is one more spatial symmetry operation which is equivalent *in two dimensions* to a 180° rotation, but which is a quite distinct operation in three dimensions. It is *inversion*. By an inversion we mean that any point at the vector displacement R from some origin [for instance, the point A in Fig. 30-9(b)] is moved to the point at $-R$.

An inversion of pattern (a) of Fig. 30-9 produces a new pattern, but an inversion of pattern (b) reproduces the same pattern. For a two-dimensional pattern (as you can see from the figure), an inversion of the pattern (b) through the point A is equivalent to a rotation of 180° about the same point. Suppose, however, we make the pattern in Fig. 30-9(b) three dimensional by imagining that the little 6's and 9's each have an “arrow” *pointing out of the page*. After an inversion in three dimensions all the arrows will be reversed, so the pattern is *not* reproduced. If we indicate the heads and tails of the arrows by dots and crosses, respectively, we can make a *three-dimensional* pattern, as in Fig. 30-9(c), which is *not* symmetric under an inversion, or we can make a pattern like the one shown in (d), which *does* have such a symmetry. Notice that it is *not* possible to imitate a three-dimensional inversion by any combination of rotations.

If we characterize the “symmetry” of a pattern—or lattice—by the kinds of symmetry operations we have been describing, it turns out that for two dimensions 17 distinct patterns are possible. We have drawn one pattern of the lowest possible

symmetry in Fig. 30-1, and one of high symmetry in Fig. 30-7. We will leave you with the game of trying to figure out all of the 17 possible patterns.

It is peculiar how few of the 17 possible patterns are used in making wallpaper and fabrics. One always sees the same three or four basic patterns. Is this because of a lack of imagination of designers, or because many of the possible patterns are not pleasing to the eye?

30-6 Symmetries in three dimensions

So far we have talked only about patterns in two dimensions. What we are really interested in, however, are patterns of atoms in three dimensions. First, it is clear that a three-dimensional crystal will have *three* primitive vectors. If we then ask about the possible symmetry operations in three dimensions, we find that there are 230 different possible symmetries! For some purposes, these 230 types can be grouped into seven classes, which are drawn in Fig. 30-10. The lattice with the least symmetry is called the *triclinic*. Its unit cell is a parallelepiped. The primitive vectors are of different lengths, and no two of the angles between them are equal. There is no possibility of any rotational or reflection symmetry. There are, however, still two possible symmetries—the unit cell is, or is not, changed by an inversion through the vertex. (By an inversion in three dimensions, we again mean that spatial displacements \mathbf{R} are replaced by $-\mathbf{R}$ —in other words, that (x, y, z) goes into $(-x, -y, -z)$). So the triclinic lattice has only two possible symmetries, unless there is some special relation among the primitive vectors. For example, if all the vectors are equal and are separated by equal angles, one has the *trigonal* lattice shown in the figure. This figure can have an additional symmetry, it may be unchanged by a rotation about the long, body diagonal.

If one of the primitive vectors, say c , is at right angles to the other two, we get a *monoclinic* unit cell. A new symmetry is possible—a rotation by 180° about c . The *hexagonal* cell is a special case in which the vectors a and b are equal and the angle between them is 60° , so that a rotation of 60° , or 120° , or 180° about the vector c repeats the same lattice (for certain internal symmetries).

If all three primitive vectors are at right angles, but of different lengths, we get the *orthorhombic* cell. The figure is symmetric for rotations of 180° about the three axes. Higher-order symmetries are possible with the *tetragonal* cell, which has all right angles and two equal primitive vectors. Finally, there is the *cubic* cell, which is the most symmetric of all.

The point of all this discussion about symmetries is that the internal symmetries of the crystals show up—sometimes in subtle ways—in the macroscopic physical properties of the crystal. For instance, a crystal will, in general, have a tensor electric polarizability. If we describe the tensor in terms of the ellipsoid of polarization, we should expect that some of the crystal symmetries should show up also in the ellipsoid. For example, a cubic crystal is symmetric with respect to a rotation of 90° about any one of three orthogonal directions. Clearly, the only ellipsoid with this property is a sphere. *A cubic crystal must be an isotropic dielectric.*

On the other hand, a tetragonal crystal has a fourfold rotational symmetry. Its ellipsoid must have two of its principal axes equal, and the third must be parallel to the axis of the crystal. Similarly, since the orthorhombic crystal has twofold rotational symmetry about three orthogonal axes, its axes must coincide with the axes of the polarization ellipsoid. In a like manner, *one* of the axes of a monoclinic crystal must be parallel to *one* of the principal axes of the ellipsoid, though we can't say anything about the other axes. Since a triclinic crystal has no rotational symmetry, the ellipsoid can have any orientation at all.

As you can see, we can make a big game of figuring out the possible symmetries and relating them to the possible physical tensors. We have considered only the polarization tensor, but things get more complicated for others—for instance, for the tensor of elasticity. There is a branch of mathematics called “group theory” that deals with such subjects, but usually you can figure out what you want with common sense.

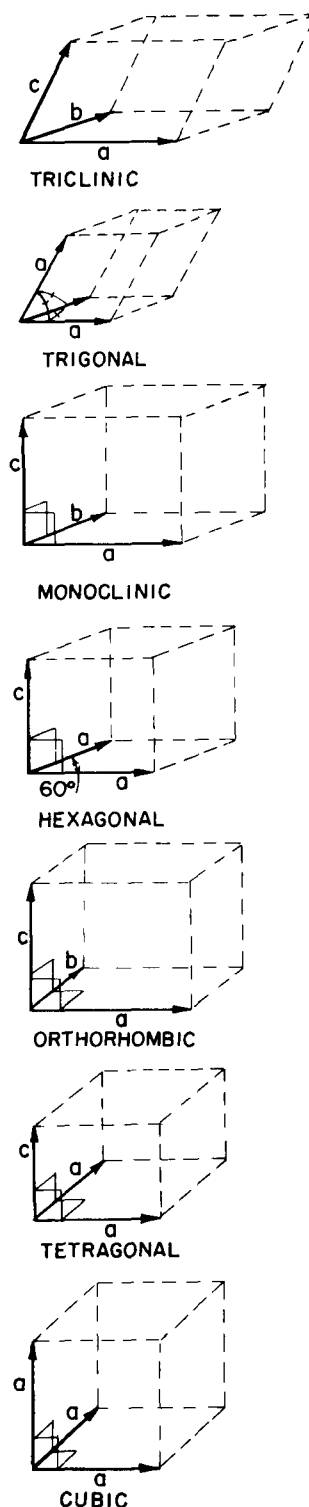


Fig. 30-10. The seven classes of crystal lattices.

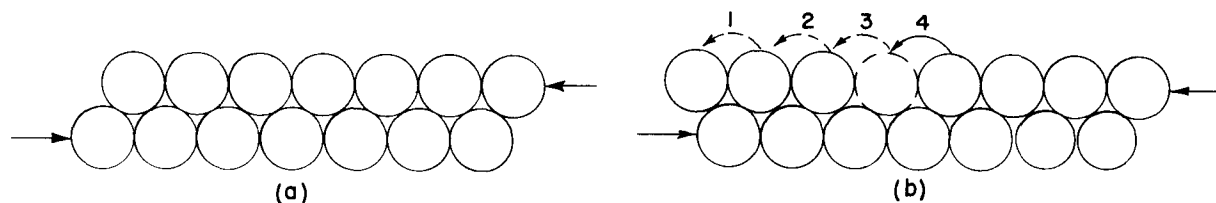


Fig. 30-11. Slippage of crystal planes.

30-7 The strength of metals

We have said that metals usually have a simple cubic crystal structure; we want now to discuss their mechanical properties—which depend on this structure. Metals are, generally speaking, very “soft,” because it is easy to slide one layer of the crystal over the next. You may think: “That’s ridiculous; metals are strong.” Not so, a *single crystal* of a metal can be distorted very easily.

Suppose we look at two layers of a crystal subjected to a shear force, as shown in the diagram of Fig. 30-11(a). You might at first think the whole layer would resist motion until the force was big enough to push the whole layer “over the hump,” so that it shifted one notch to the left. Although slipping does occur along a plane, it doesn’t happen that way (If it did, you would calculate that the metal is much stronger than it really is.) What happens is more like one atom going at a time; first the atom on the left makes its jump, then the next, and so on, as indicated in Fig. 30-11(b). In effect it is the vacant space between two atoms that quickly travels to the right, with the net result that the whole second layer has moved over one atomic spacing. The slipping goes this way because it takes much less energy to lift one atom at a time over the hump than to lift a whole row. Once the force is enough to start the process, it goes the rest of the way very fast.

It turns out that in a real crystal, slipping will occur repeatedly at one plane, then will stop there and start at some other plane. The details of why it starts and stops are quite mysterious. It is, in fact, quite strange that successive regions of slip are often fairly evenly spaced. Figure 30-12 shows a photograph of a tiny, thin copper crystal that has been stretched. You can see the various planes where slipping has occurred.

Fig. 30-12. A photograph of a small crystal of copper after stretching. [Courtesy of S. S. Brenner, Senior Scientist, United States Steel Research Center, Monroeville, Pa.]

The sudden slipping of individual crystal planes is quite apparent if you take a piece of tin wire that has large crystals in it and stretch it while holding it next to your ear. You can hear a rush of “ticks” as the planes snap to their new positions, one after the other.

The problem of having a “missing” atom in one row is somewhat more difficult than it might appear from Fig. 30-11. When there are more layers, the situation must be something like that shown in Fig. 30-13. Such an imperfection in a crystal is called a *dislocation*. It is presumed that such dislocations are either present when the crystal was formed or are generated at some notch or crack at the surface. Once they are produced, they can move relatively freely through the crystal. The gross distortions result from the motions of many of such dislocations.

Dislocations can move freely—that is, they require little extra energy—so long as the rest of the crystal has a perfect lattice. But they may get “stuck” if they encounter some other kind of imperfection in the crystal. If it takes a lot of energy for them to pass the imperfection, they will be stopped. This is precisely the mechanism that gives strength to *imperfect* metal crystals. Pure iron crystals are quite soft, but a small concentration of impurity atoms may cause enough imperfections to effectively immobilize the dislocations. As you know, steel, which is primarily iron, is very hard. To make steel, a small amount of carbon is dissolved in the iron melt; if the melt is cooled rapidly, the carbon precipitates out in little grains, making many microscopic distortions in the lattice. The dislocations can no longer move about, and the metal is hard.

Pure copper is very soft, but can be “work-hardened.” This is done by hammering on it or bending it back and forth. In this case, many new dislocations of various kinds are made which interfere with one another, cutting down their

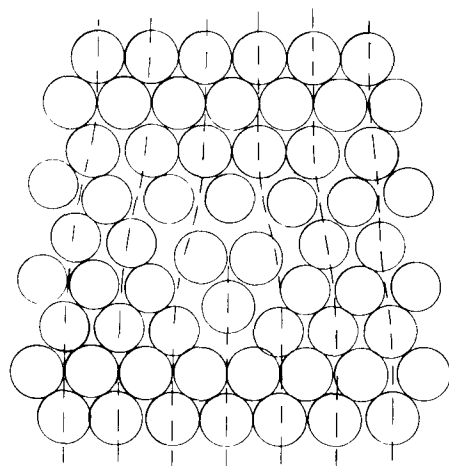


Fig. 30-13. A dislocation in a crystal.

mobility. Perhaps you've seen the trick of taking a bar of "dead soft" copper and gently bending it around someone's wrist as a bracelet. In the process, it becomes work-hardened and cannot easily be unbent again! A work-hardened metal like copper can be made soft again by annealing at a high temperature. The thermal motion of the atoms "irons out" the dislocations and makes large single crystals again. We have, so far, described only the so-called *slip* dislocation. There are many other kinds, one of which is the *screw* dislocation shown in Fig. 30-14. Such dislocations often play an important part in crystal growth.

30-8 Dislocations and crystal growth

One of the great puzzles for a long time was how crystals can possibly grow. We have described how it is that each atom might, by repeated testing, determine whether it was better to be in the crystal or not. But that means that each atom must find a place of low energy. However, an atom put on a new surface is only bound by one or two bonds from below, and doesn't have the same energy it would have if it were placed in a corner, where it would have atoms on three sides. Suppose we imagine a growing crystal as a stack of blocks, as shown in Fig. 30-15. If we try a new block at, say, position *A*, it will have only one of the six neighbors it should ultimately get. With so many bonds lacking, its energy is not very low. It would be better off at position *B*, where it already has one-half of its quota of bonds. Crystals do indeed grow by attaching new atoms at places like *B*.

What happens, though, when that line is finished? To start a new line, an atom must come to rest with only two sides attached, and that is again not very likely. Even if it did, what would happen when the layer was finished? How could a new layer get started? One answer is that the crystal prefers to grow at a dislocation, for instance around a screw dislocation like the one shown in Fig. 30-14. As blocks are added to this crystal, there is always some place where there are three available bonds. The crystal prefers, therefore, to grow with a dislocation built in. Such a spiral pattern of growth is shown in Fig. 30-16, which is a photograph of a single crystal of paraffin.

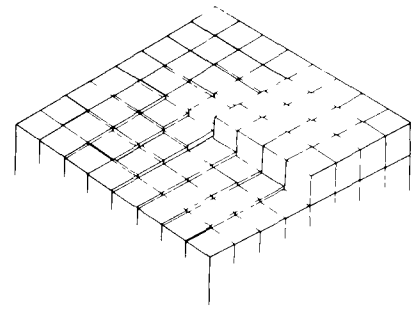


Fig. 30-14. A screw dislocation. [From Charles Kittel, *Introduction to Solid State Physics*, John Wiley and Sons, Inc., New York, 2nd ed., 1956.]

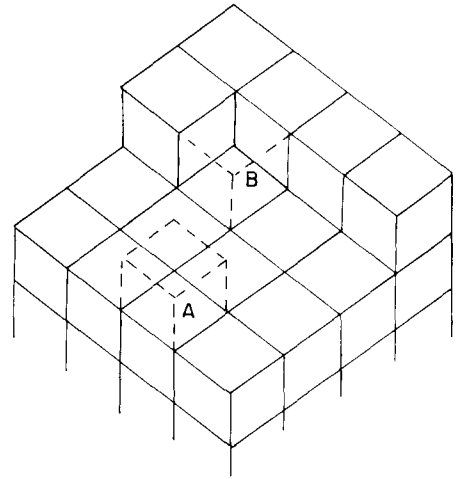


Fig. 30-15. Crystal growth.

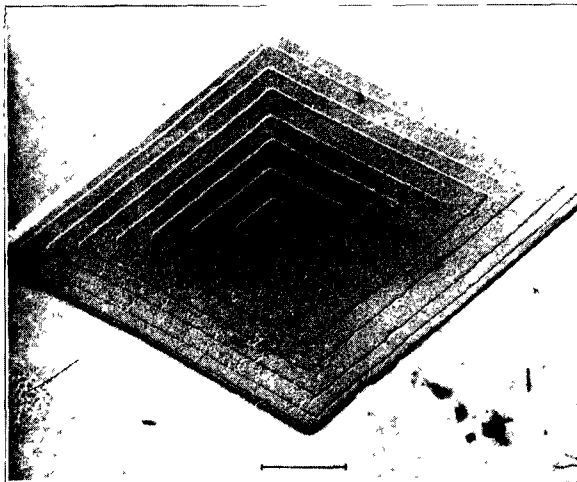


Fig. 30-16. A paraffin crystal which has grown around a screw dislocation. [From Charles Kittel, *Introduction to Solid State Physics*, John Wiley and Sons, Inc., New York, 2nd ed., 1956.]

30-9 The Bragg-Nye crystal model

We cannot, of course, see what goes on with the individual atoms in a crystal. Also, as you realize by now, there are many complicated phenomena that are not easy to treat quantitatively. Sir Lawrence Bragg and J. F. Nye have devised a scheme for making a model of a metallic crystal which shows in a striking way many of the phenomena that are believed to occur in a real metal. In the following pages we have reproduced their original article, which describes their method and shows some of the results they obtained with it. (The article is reprinted from the *Proceedings of the Royal Society of London*, Vol. 190, September 1947, pp. 474-481—with the permission of the authors and of the Royal Society.)

A dynamical model of a crystal structure

BY SIR LAWRENCE BRAGG, F.R.S. AND J. F. NYE

Cavendish Laboratory, University of Cambridge

(Received 9 January 1947—Read 19 June 1947)

[Plates 8 to 21]

The crystal structure of a metal is represented by an assemblage of bubbles, a millimetre or less in diameter, floating on the surface of a soap solution. The bubbles are blown from a fine pipette beneath the surface with a constant air pressure, and are remarkably uniform in size. They are held together by surface tension, either in a single layer on the surface or in a three-dimensional mass. An assemblage may contain hundreds of thousands of bubbles and persists for an hour or more. The assemblages show structures which have been supposed to exist in metals, and simulate effects which have been observed, such as grain boundaries, dislocations and other types of fault, slip, recrystallization, annealing, and strains due to 'foreign' atoms.

1 THE BUBBLE MODEL

Models of crystal structure have been described from time to time in which the atoms are represented by small floating or suspended magnets, or by circular disks floating on a water surface and held together by the forces of capillary attraction. These models have certain disadvantages, for instance, in the case of floating objects in contact, frictional forces impede their free relative movement. A more serious disadvantage is that the number of components is limited for a large number of components is required in order to approach the state of affairs in a real crystal. The present paper describes the behaviour of a model in which the atoms are represented by small bubbles from 2.0 to 0.1 mm in diameter floating on the surface of a soap solution. These small bubbles are sufficiently persistent for experiments lasting an hour or more, they slide past each other without friction, and they can be produced in large numbers. Some of the illustrations in this paper were taken from assemblages of bubbles numbering 100,000 or more. The model most nearly represents the behaviour of a metal structure, because the bubbles are of one type only and are held together by a general capillary attraction which represents the binding force of the free electrons in the metal. A brief description of the model has been given in the *Journal of Scientific Instruments* (Bragg 1942b).

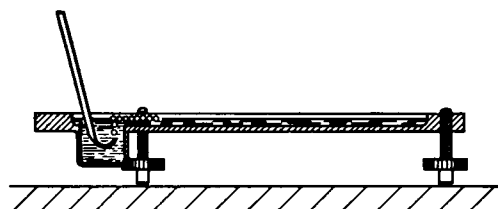


FIGURE 1. Apparatus for producing rafts of bubbles.

2 METHOD OF FORMATION

The bubbles are blown from a fine orifice, beneath the surface of a soap solution. We have had the best results with a solution the formula of which was given to us by Mr Green of the Royal Institution. 15.2 c.c. of oleic acid (pure redistilled) is well shaken in 50 c.c. of distilled water. This is mixed thoroughly with 73 c.c. of 10% solution of tri-ethanolamine and the mixture made up to 200 c.c. To this is added 164 c.c. of pure glycerine. It is left to stand and the clear liquid is drawn off from below. In some experiments this was diluted in three times its volume of water to reduce viscosity. The orifice of the jet is about 5 mm. below the surface. A constant air pressure of 50 to 200 cm. of water is supplied by means of two Winchester flasks. Normally the bubbles are remarkably uniform in size. Occasionally they issue in an irregular manner, but this can be corrected by a change of jet or of pressure. Unwanted bubbles can easily be destroyed by playing a small flame over the surface. Figure 1 shows the apparatus. We have found it of advantage to blacken the bottom of the vessel, because details of structure, such as grain boundaries and dislocations, then show up more clearly.

Figure 2, plate 8, shows a portion of a 'raft' or two-dimensional crystal of bubbles. Its regularity can be judged by looking at the figure in a glancing direction. The size of the bubbles varies with the aperture, but does not appear to vary to any marked degree with the pressure or the depth of the orifice beneath the surface. The main effect of increasing the pressure is to increase the rate of issue of the bubbles. As an example, a thick-walled jet of 49 μ bore with a pressure of 100 cm. produced bubbles of 1.2 mm. in diameter. A thin-walled jet of 27 μ diameter and a pressure of 180 cm. produced bubbles of 0.6 mm. diameter. It is convenient to refer to bubbles of 2.0 to 1.0 mm. diameter as 'large' bubbles, those from 0.8 to 0.6 mm. diameter as 'medium' bubbles, and those from 0.3 to 0.1 mm. diameter as 'small' bubbles, since their behaviour varies with their size.

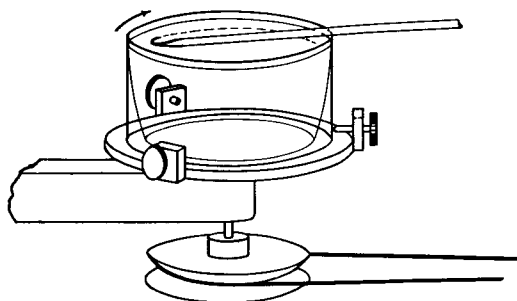


FIGURE 3. Apparatus for producing bubbles of small size.

With this apparatus we have not found it possible to reduce the size of the jet and so produce bubbles of smaller diameter than 0.6 mm. As it was desired to experiment with very small bubbles, we had recourse to placing the soap solution in a rotating vessel and introducing a fine jet as nearly as possible parallel to a stream line. The bubbles are swept away as they form, and under steady conditions are reasonably uniform. They issue at a rate of one thousand or more per second, giving a high-pitched note. The soap solution mounts up in a steep wall around the perimeter of the vessel while it is rotating but carries back most of the bubbles with it when rotation ceases. With this device, illustrated in figure 3, bubbles down to 0.12 mm. in diameter can be obtained. As an example, an orifice 38 μ across in a thin-walled jet, with a pressure of 190 cm. of water, and a speed of the fluid of 180 cm./sec. past the orifice, produced bubbles of 0.14 mm. diameter. In this case a dish of diameter 9.5 cm. and speed of 6 rev./sec. was used. Figure 4, plate 8, is an enlarged picture of these 'small' bubbles and shows their degree of regularity, the pattern is not as perfect with a rotating as with a stationary vessel, the rows being seen to be slightly irregular when viewed in a glancing direction.

These two-dimensional crystals show structures which have been supposed to exist in metals, and simulate effects which have been observed, such as grain boundaries, dislocations and other types of fault, slip, recrystallization, annealing, and strains due to 'foreign' atoms.

3 GRAIN BOUNDARIES

Figures 5a, 5b and 5c, plates 9 and 10, show typical grain boundaries for bubbles of 1.87, 0.76 and 0.30 mm. diameter respectively. The width of the disturbed area at the boundary, where the bubbles have an irregular distribution, is in general greater the smaller the bubbles. In figure 5a, which shows portions of several adjacent grains, bubbles at a boundary between two grains adhere definitely to one crystalline arrangement or the other. In figure 5c there is a marked 'Beilby layer' between the two grains. The small bubbles, as will be seen, have a greater rigidity than the large ones, and this appears to give rise to more irregularity at the interface.

Separate grains show up distinctly when photographs of polycrystalline rafts such as figures 5a to 5c, plates 9 and 10, and figures 12a to 12e, plates 14 to 16, are viewed obliquely. With suitable lighting, the floating raft of bubbles itself when viewed obliquely resembles a polished and etched metal in a remarkable way.

It often happens that some 'impurity atoms', or bubbles which are markedly larger or smaller than the average, are found in a polycrystalline raft, and when this is so a large proportion of them are situated at the grain boundaries. It would be incorrect to say that the irregular bubbles make their way to the boundaries, it is a defect of the model that no diffusion of bubbles through the structure can take place, mutual adjustments of neighbours alone being possible. It appears that the boundaries tend to readjust themselves by the growth of one crystal at the expense of another till they pass through the irregular atoms.

4 DISLOCATIONS

When a single crystal or polycrystalline raft is compressed, extended, or otherwise deformed it exhibits a behaviour very similar to that which has been pictured for metals subjected to strain. Up to a certain limit the model is within its elastic range. Beyond that point it yields by slip along one of the three equally inclined directions of closely packed rows. Slip takes place by the bubbles in one row moving forward over those in the next row by an amount equal to the distance between neighbours. It is very interesting to watch this process taking place. The movement is not simultaneous along the whole row but begins at one end with the appearance of a 'dislocation', where there is locally one more bubble in the rows on one side of the slip line as compared with those on the other. This dislocation then runs along the slip line from one side of the crystal to the other, the final result being a slip by one 'inter-atomic' distance. Such a process has been invoked by Orowan, by Polanyi and by Taylor to explain the small forces required to produce plastic gliding in metal structures. The theory put forward by Taylor (1934) to explain the mechanism of plastic deformation of crystals considers the mutual action and equilibrium of such dislocations. The bubbles afford a very striking picture of what has been supposed to take place in the metal. Sometimes

the dislocations run along quite slowly, taking a matter of seconds to cross a crystal, stationary dislocations also are to be seen in crystals which are not homogeneously strained. They appear as short black lines, and can be seen in the series of photographs, figures 12a to 12e, plates 14 to 16. When a polycrystalline raft is compressed, these dark lines are seen to be dashing about in all directions across the crystals.

Figures 6a, 6b and 6c, plates 10 and 11, show examples of dislocations. In figure 6a, where the diameter of the bubbles is 1.9 mm, the dislocation is very local, extending over about six bubbles. In figure 6b (diameter 0.76 mm) it extends over twelve bubbles, and in figure 6c (diameter 0.30 mm) its influence can be traced for a length of about fifty bubbles. The greater rigidity of the small bubbles leads to longer dislocations. The study of any mass of bubbles shows, however, that there is not a standard length of dislocation for each size. The length depends upon the nature of the strain in the crystal. A boundary between two crystals with corresponding axes at approximately 30° (the maximum angle which can occur) may be regarded as a series of dislocations in alternate rows, and in this case the dislocations are very short. As the angle between the neighbouring crystals decreases, the dislocations occur at wider intervals and at the same time become longer, till one finally has single dislocations in a large body of perfect structure as shown in figures 6a, 6b and 6c.

Figure 7, plate 11, shows three parallel dislocations. If we call them positive and negative (following Taylor) they are positive, negative, positive, reading from left to right. The strip between the last two has three bubbles in excess, as can be seen by looking along the rows in a horizontal direction. Figure 8, plate 12, shows a dislocation projecting from a grain boundary, an effect often observed.

Figure 9, plate 12, shows a place where two bubbles take the place of one. This may be regarded as a limiting case of positive and negative dislocations on neighbouring rows, with the compressive sides of the dislocations facing each other. The contrary case would lead to a hole in the structure, one bubble being missing at the point where the dislocations met.

5. OTHER TYPES OF FAULT

Figure 10, plate 12, shows a narrow strip between two crystals of parallel orientation, the strip being crossed by a number of fault lines where the bubbles are not in close packing. It is in such places as these that recrystallization may be expected. The boundaries approach and the strip is absorbed into a wider area of perfect crystal.

Figures 11a to 11g, plates 13 and 14 are examples of arrangements which frequently appear in places where there is a local deficiency of bubbles. While a dislocation is seen as a dark stripe in a general view, these structures show up in the shape of the letter V or as triangles. A typical V structure is seen in figure 11a. When the model is being distorted, a V structure is formed by two dislocations meeting at an inclination of 60°, it is destroyed by the dislocations continuing along their paths. Figure 11b shows a small triangle, which also embodies a dislocation, for it will be noticed that the rows below the fault have one more bubble than those below. If a mild amount of 'thermal movement' is imposed by gentle agitation of one side of the crystal, such faulty places disappear and a perfect structure is formed.

Here and there in the crystals there is a blank space where a bubble is missing, showing as a black dot in a general view. Examples occur in figure 11g. Such a gap cannot be closed by a local readjustment, since filling the hole causes another to appear. Such holes both appear and disappear when the crystal is 'cold-worked'.

These structures in the model suggest that similar local faults may exist in an actual metal. They may play a part in processes such as diffusion or the order-disorder change by reducing energy barriers in their neighbourhood, and act as nuclei for crystallization in an allotropic change.

6. RECRYSTALLIZATION AND ANNEALING

Figures 12a to 12e, plates 14 to 16, show the same raft of bubbles at successive times. A raft covering the surface of the solution was given a vigorous stirring with a glass rake, and then left to adjust itself. Figure 12a shows its aspect about 1 sec. after stirring has ceased. The raft is broken into a number of small 'crystallites', these are in a high state of non-homogeneous strain as is shown by the numerous dislocations and other faults. The following photograph (figure 12b) shows the same raft 32 sec. later. The small grains have coalesced to form larger grains, and much of the strain has disappeared in the process. Recrystallization takes place right through the series, the last three photographs of which show the appearance of the raft 2, 14 and 25 min. after the initial stirring. It is not possible to follow the rearrangement for much longer times, because the bubbles shrink after long standing, apparently due to the diffusion of air through their walls, and they also become thin and tend to burst. No agitation was given to the model during this process. An ever slower process of rearrangement goes on, the movement of the bubbles in one part of the raft setting up strains which activate a rearrangement in a neighbouring part, and that in its turn still another.

A number of interesting points are to be seen in this series. Note the three small grains at the points indicated by the co-ordinates AA, BB, CC. A persists, though

changed in form, throughout the whole series. B is still present after 14 min., but has disappeared in 25 min., leaving behind it four dislocations marking internal strain in the grain. Grain C shrinks and finally disappears in figure 12d, leaving a hole and a V which has disappeared in figure 12e. At the same time the ill-defined boundary in figure 12d at DD has become a definite one in figure 12e. Note also the straightening out of the grain boundary in the neighbourhood of EE in figures 12b to 12e. Dislocations of various lengths can be seen, marking all stages between a slight warping of the structure and a definite boundary. Holes where bubbles are missing show up as black dots. Some of these holes are formed or filled up by movements of dislocations, but others represent places where a bubble has burst. Many examples of V's and some of triangles can be seen. Other interesting points will be apparent from a study of this series of photographs.

Figures 13a, 13b and 13c, plate 17, show a portion of a raft 1 sec., 4 sec. and 4 min. after the stirring process, and is interesting as showing two successive stages in the relaxation towards a more perfect arrangement. The changes show up well when one looks in a glancing direction across the page. The arrangement is very broken in figure 13a. In figure 13b the bubbles have grouped themselves in rows, but the curvature of these rows indicates a high degree of internal strain. In figure 13c this strain has been relieved by the formation of a new boundary at A-A, the rows on either side now being straight. It would appear that the energy of this strained crystal is greater than that of the intercrystalline boundary. We are indebted to Messrs Kodak for the photographs of figure 13, which were taken when the cinematograph film referred to below was produced.

7. EFFECT OF IMPURITY ATOM

Figure 14, plate 18, shows the widespread effect of a bubble which is of the wrong size. If this figure is compared with the perfect rafts shown in figures 2 and 4, plate 8, it will be seen that three bubbles, one larger and two smaller than normal, disturb the regularity of the rows over the whole of the figure. As has been mentioned above, bubbles of the wrong size are generally found in the grain boundaries, where holes of irregular size occur which can accommodate them.

8. MECHANICAL PROPERTIES OF THE TWO-DIMENSIONAL MODEL

The mechanical properties of a two-dimensional perfect raft have been described in the paper referred to above (Bragg 1942b). The raft lies between two parallel springs dipping horizontally in the surface of the soap solution. The pitch of the springs is adjusted to fit the spacing of the rows of bubbles, which then adhere firmly to them. One spring can be translated parallel to itself by a micrometer screw, and the other is supported by two thin vertical glass fibres. The shearing stress can be measured by noting the deflexion of the glass fibres. When subjected to a shearing strain, the raft obeys Hooke's law of elasticity up to the point where the elastic limit is reached. It then slips along some intermediate row by an amount equal to the width of one bubble. The elastic shear and slip can be repeated several times. The elastic limit is approximately reached when one side of the raft has been sheared by an amount equal to a bubble width past the other side. This feature supports the basic assumption made by one of us in the calculation of the elastic limit of a metal (Bragg 1942a), in which it is supposed that each crystallite in a cold-worked metal only yields when the strain in it has reached such a value that energy is released by the slip.

A calculation has been made by M. M. Nicolson of the forces between the bubbles, and will be published shortly. It shows two interesting points. The curve for the variation of potential energy with distance between centres is very similar to those which have been plotted for atoms. It has a minimum for a distance between centres slightly less than a free bubble diameter, and rises sharply for smaller distances. Further, the rise is extremely sharp for bubbles of 0.1 mm. diameter but much less so for bubbles of 1 mm. diameter, thus confirming the impression given by the model that the small bubbles behave as if they were much more rigid than the large ones.

9. THREE-DIMENSIONAL ASSEMBLAGES

If the bubbles are allowed to accumulate in multiple layers on the surface, they form a mass of three-dimensional 'crystals' with one of the arrangements of closest packing. Figure 15, plate 18, shows an oblique view of such a mass, its resemblance to a polished and etched metal surface is noticeable. In figure 16, plate 20, a similar mass is seen viewed normally. Parts of the structure are definitely in cubic closest packing, the outer surface being the (111) face or (100) face. Figure 17a, plate 19, shows a (111) face. The outlines of the three bubbles on which each upper bubble rests can be clearly seen, and the next layer of these bubbles is faintly visible in a position not beneath the uppermost layer, showing that the packing of the (111) planes has the well-known cubic succession. Figure 17b, plate 19, shows a (100) face with each bubble resting on four others. The cubic axes are of course inclined at 45° to the close-packed rows of the surface layer. Figure 17c, plate 19, shows a twin in the cubic structure across the face (111). The uppermost faces are (111) and (100), and they make a small angle with each other, though this is not apparent in the figure, it shows up in an oblique view. Figure 17d, plate 19, appears to show both the cubic and hexagonal succession of closely packed planes, but it is difficult to verify whether the left-hand side follows the true hexagonal close-packed struc-

ture because it is not certain that the assemblage had a depth of more than two layers at this point. Many instances of twins, and of intercrystalline boundaries, can be seen in figure 16, plate 20.

Figure 18, plate 21, shows several dislocations in a three-dimensional structure subjected to a bending strain.

10 DEMONSTRATION OF THE MODEL

With the co-operation of Messrs Kodak, a 16 mm cinematograph film has been made of the movements of the dislocations and grain boundaries when single crystal and polycrystalline rafts are sheared, compressed, or extended. Moreover, if the soap solution is placed in a glass vessel with a flat bottom, the model lends itself to projection on a large scale by transmitted light. Since a certain depth is required for producing the bubbles, and the solution is rather opaque, it is desirable to make the projection through a glass block resting on the bottom of the vessel and just submerged beneath the surface.

In conclusion, we wish to express our thanks to Mr C. E. Harrold, of King's College, Cambridge, who made for us some of the pipettes which were used to produce the bubbles.

REFERENCES

- Bragg, W. L. 1942a *Nature*, **149**, 511.
Bragg, W. L. 1942b *J. Sci. Instrum.* **19**, 148.
Taylor, G. I. 1934 *Proc. Roy. Soc. A*, **145**, 362.

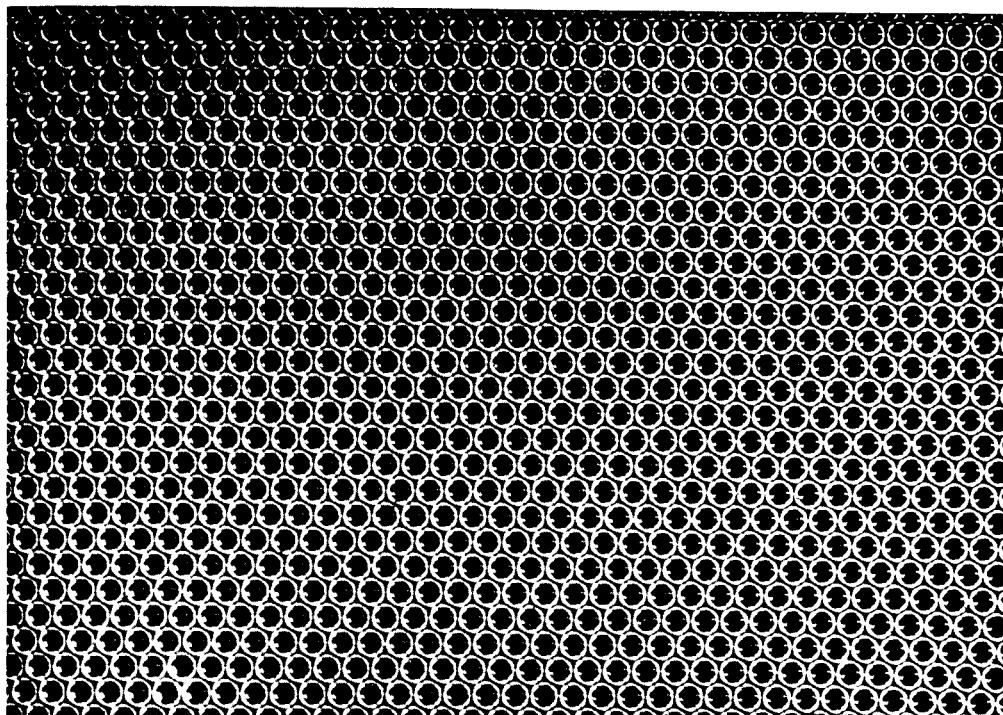


FIGURE 2. Perfect crystalline raft of bubbles. Diameter 1.41 mm.

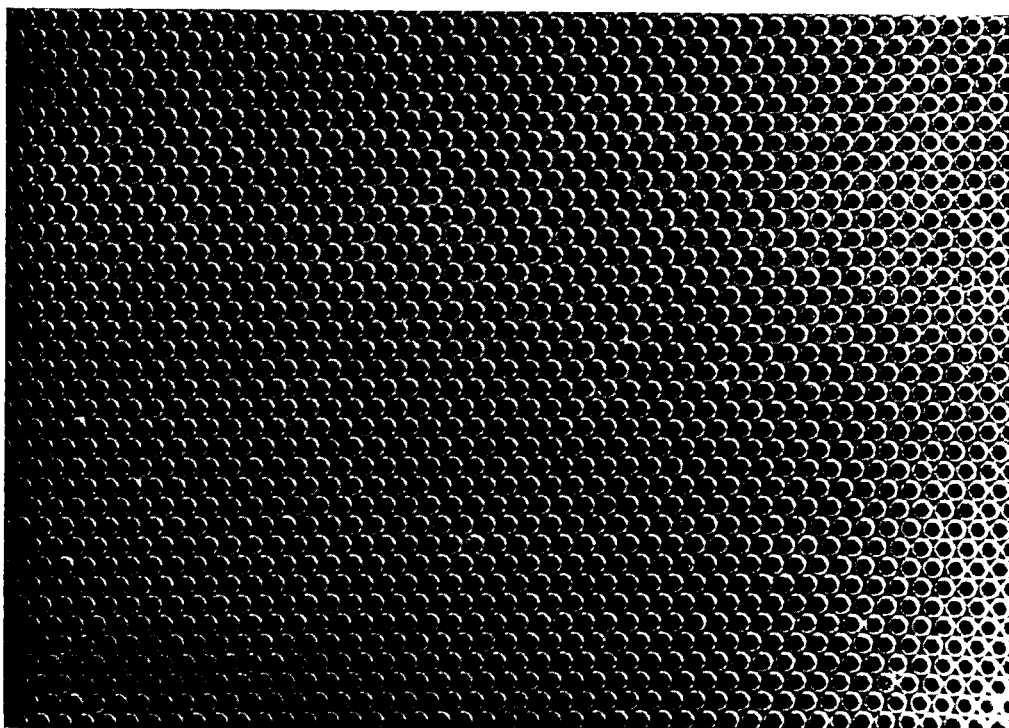


FIGURE 4. Perfect crystalline raft of bubbles. Diameter 0.30 mm.

Grain boundaries

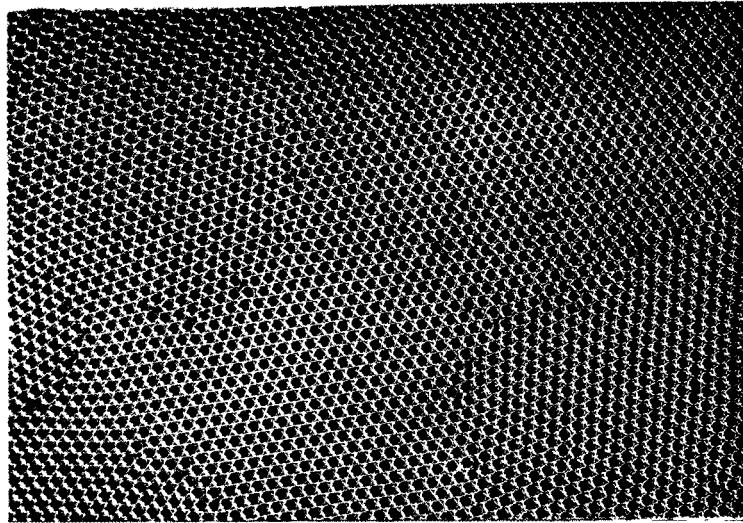


FIGURE 5a. Diameter 1.87 mm.

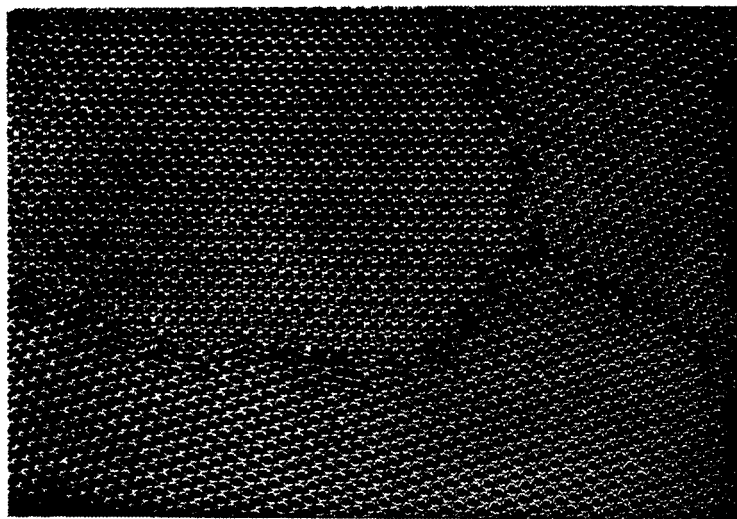


FIGURE 5b. Diameter 0.76 mm.

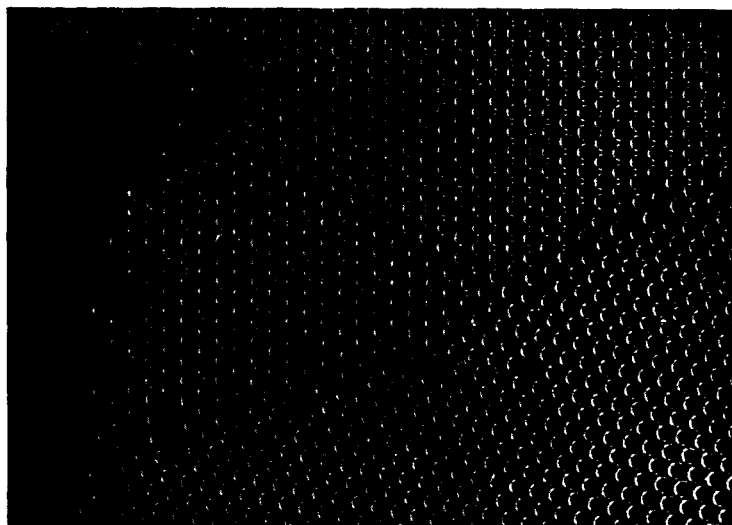


FIGURE 5*c*. A grain boundary. Diameter 0.30 mm.

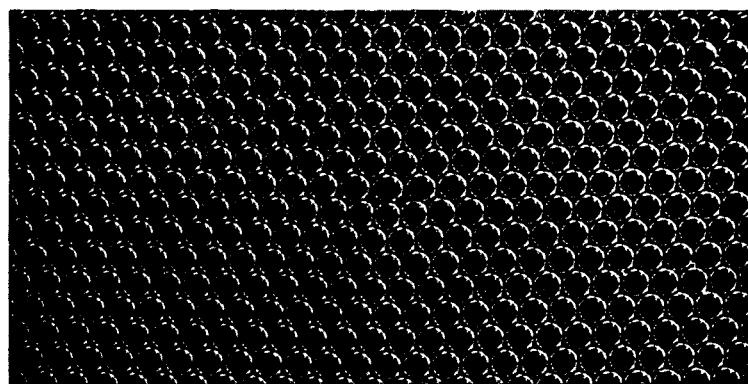


FIGURE 6*a*. A dislocation. Diameter 1.9 mm.

Dislocations

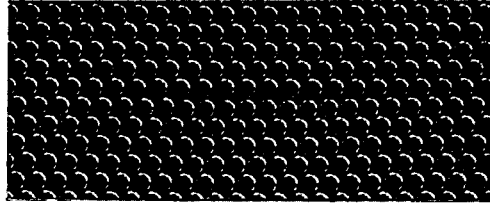


FIGURE 6*b*. Diameter 0.76 mm.

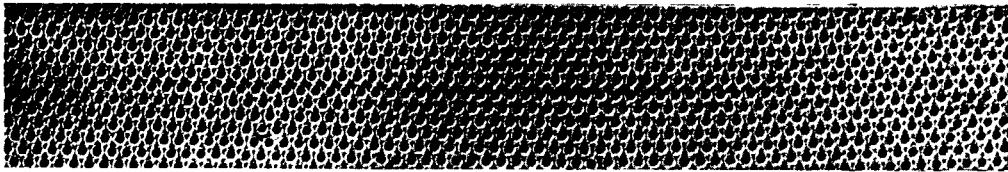


FIGURE 6*c*. Diameter 0.30 mm.

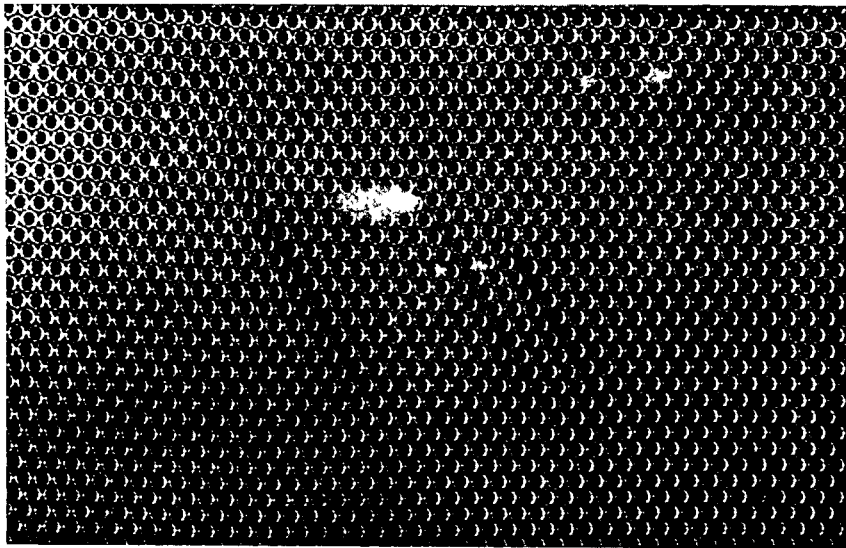


FIGURE 7. Parallel dislocations. Diameter 0.76 mm.

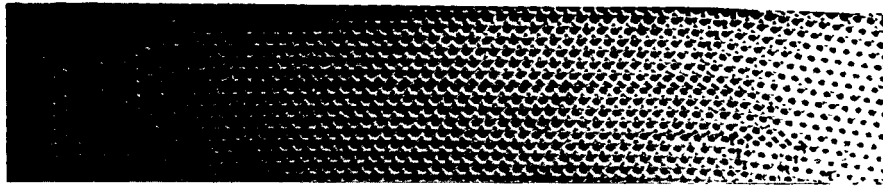


FIGURE 8. Dislocation projecting from a grain boundary. Diameter 0.30 mm.

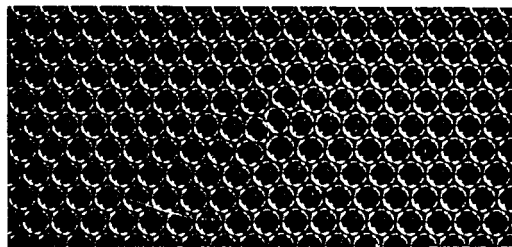


FIGURE 9. Dislocations in adjacent rows. Diameter 1.9 mm.

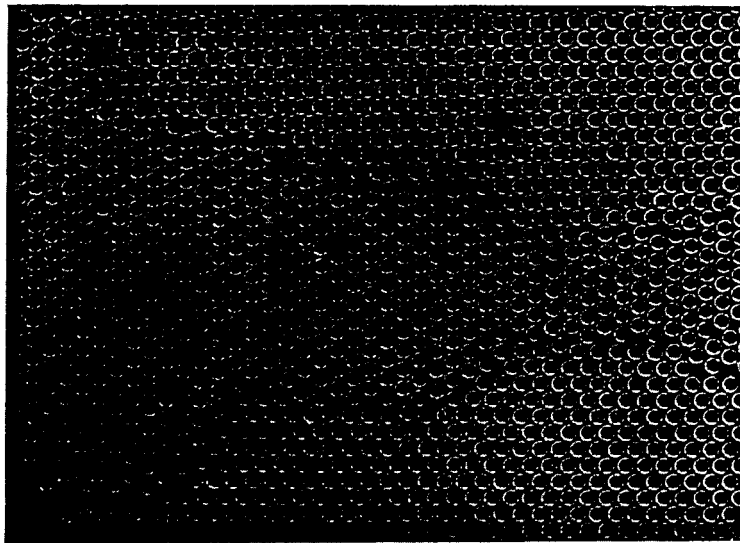
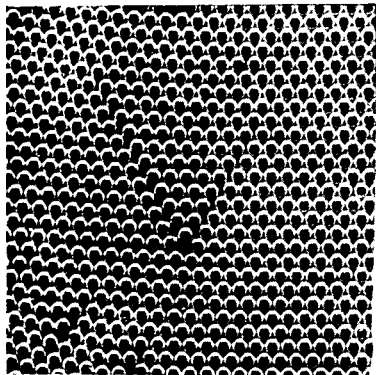
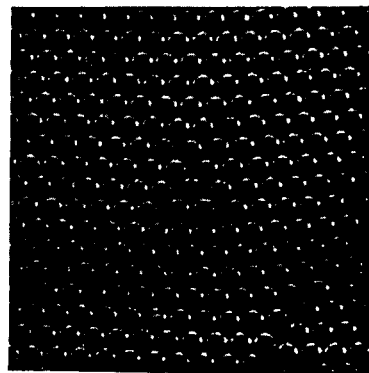


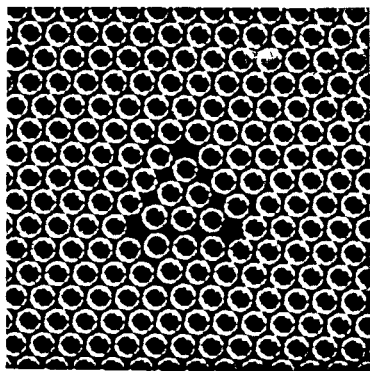
FIGURE 10. Series of fault lines between two areas of parallel orientation. Diameter 0.30 mm.



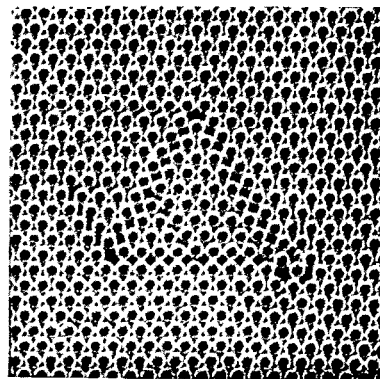
Diameter 0.68 mm.
a



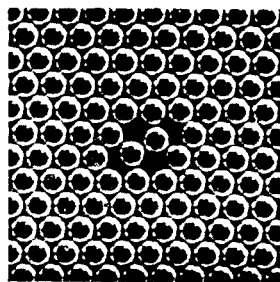
Diameter 0.68 mm.
b



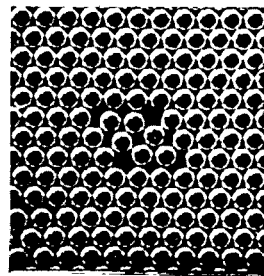
Diameter 0.6 mm.
c



Diameter 0.30 mm.
d

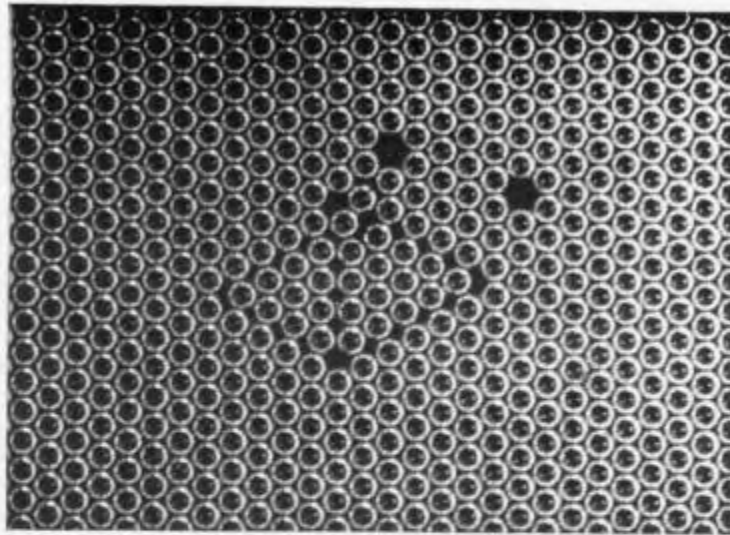


Diameter 0.6 mm
e

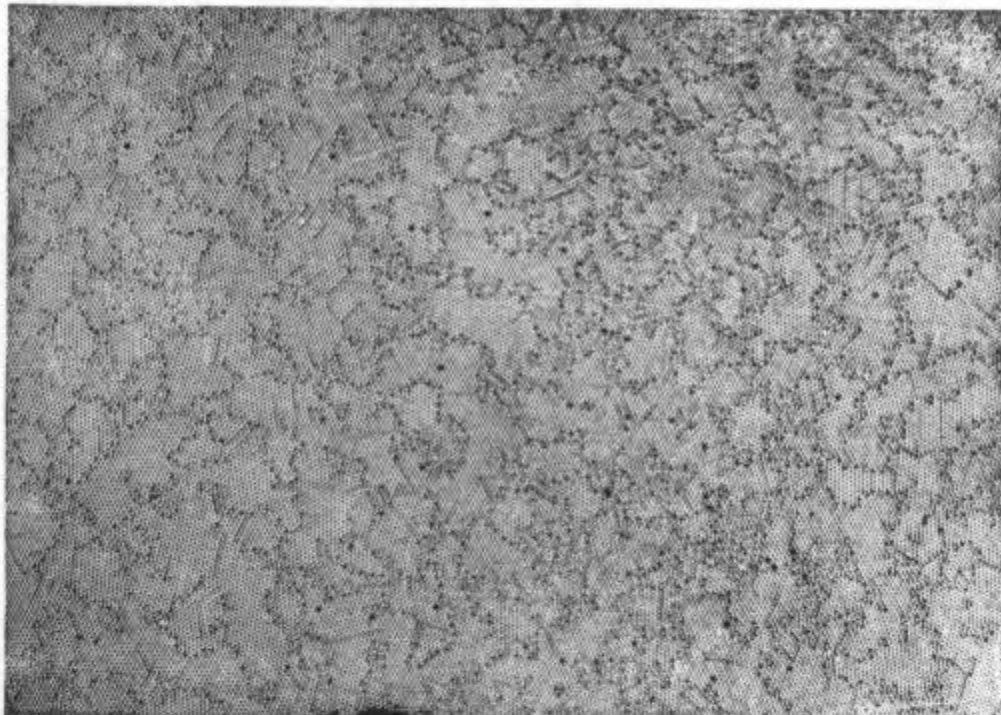


Diameter 0.6 mm
f

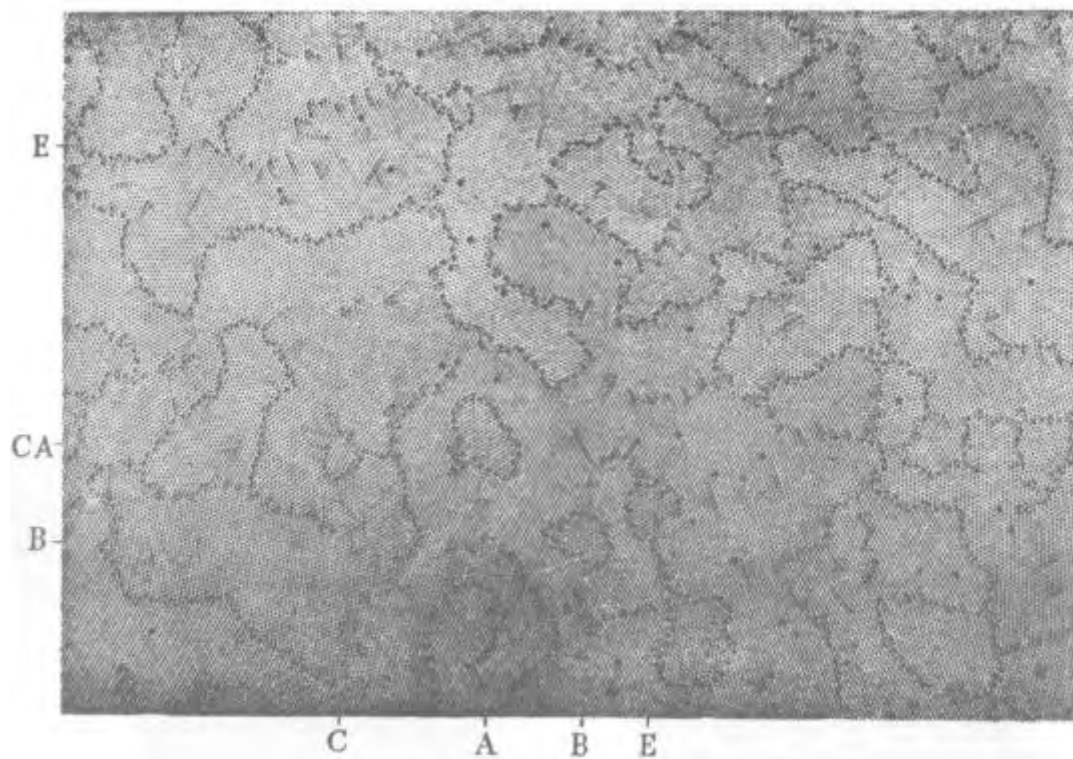
FIGURE 11. Types of fault.



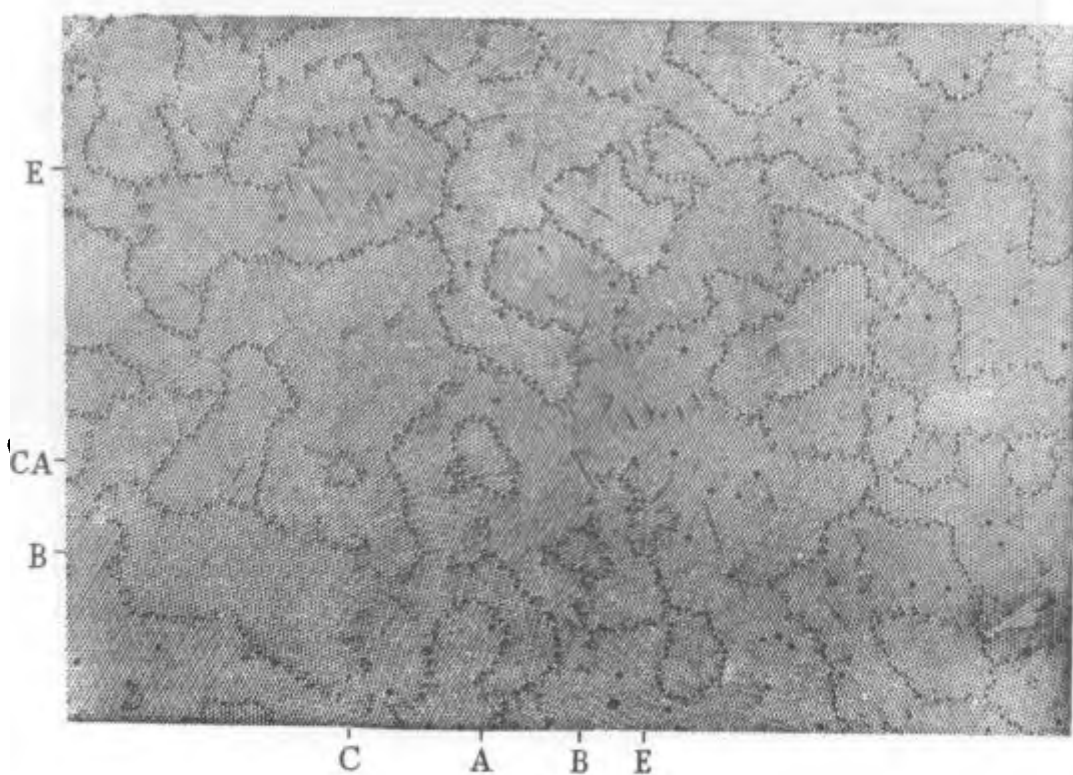
g
FIGURE 11. Types of fault. Diameter 0.68 mm.



a. Immediately after stirring;
FIGURE 12. Recrystallization. Diameter 0.60 mm.

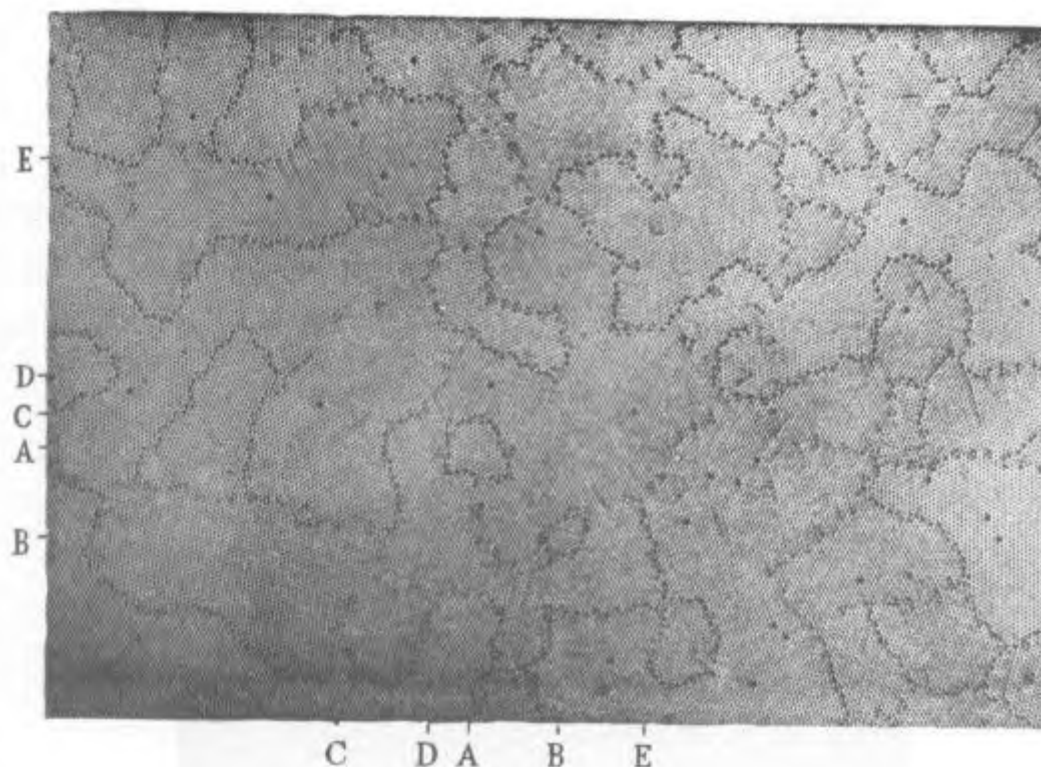


b. After 33 sec.

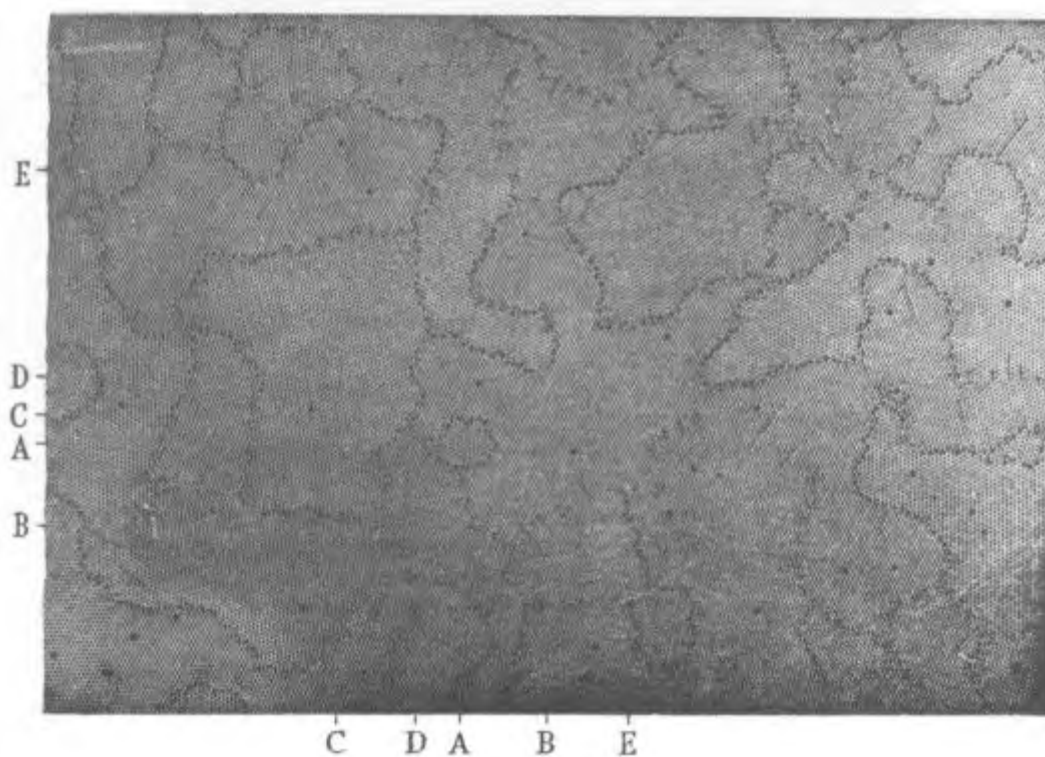


c. After 2 min.

FIGURE 12

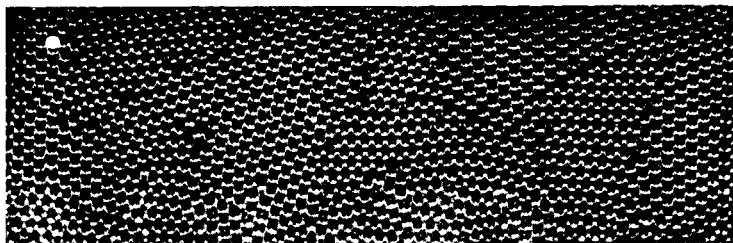


d. After 14 min.

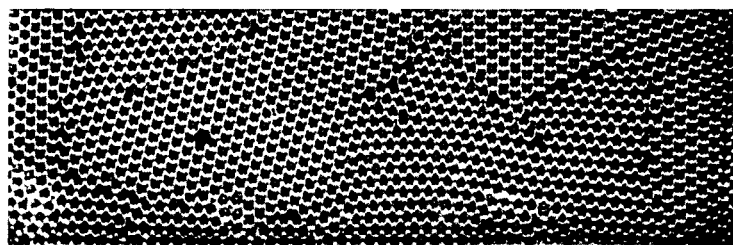


e. After 25 min.

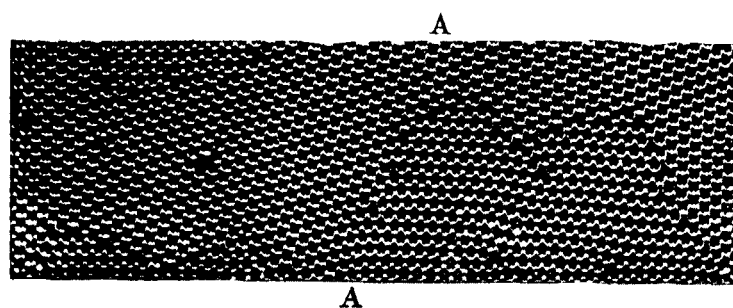
FIGURE 12



a After 1 sec.



b After 4 sec.



c After 4 min.

FIGURE 13. Two stages of recrystallization. Diameter 1.64 mm.

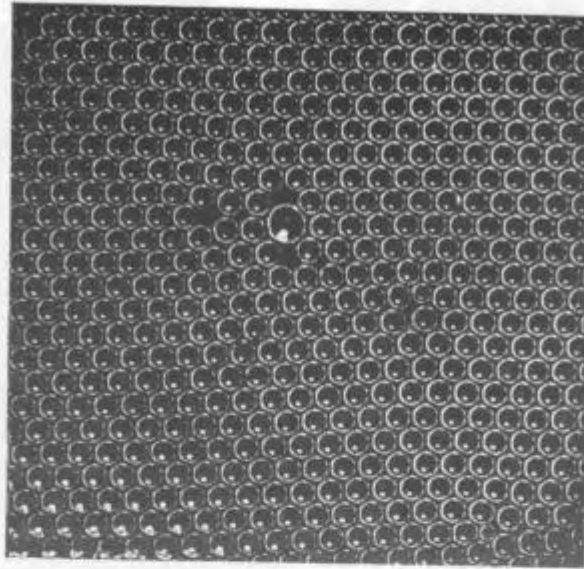
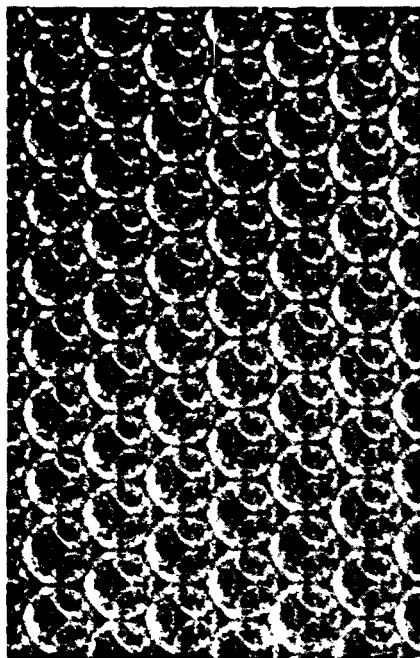


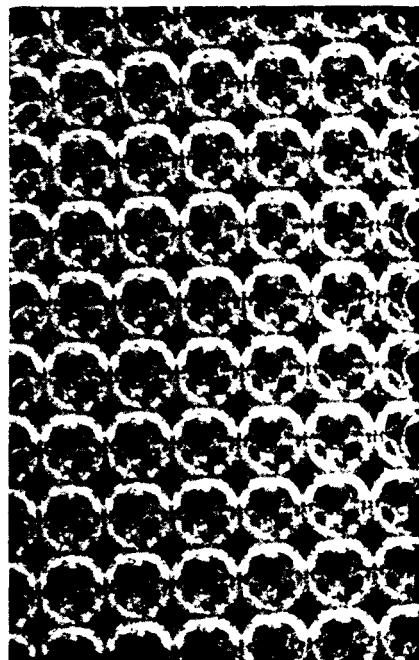
FIGURE 14. Effect of atoms of impurity. Diameter of uniform bubbles about 1.3 μ m.



FIGURE 15. Oblique view of three-dimensional raft.

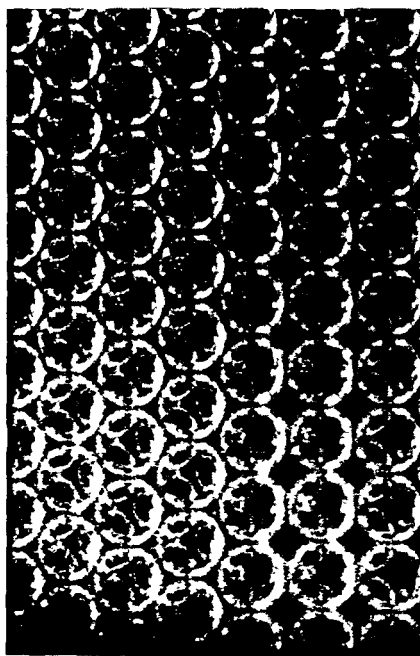


a. (111) face

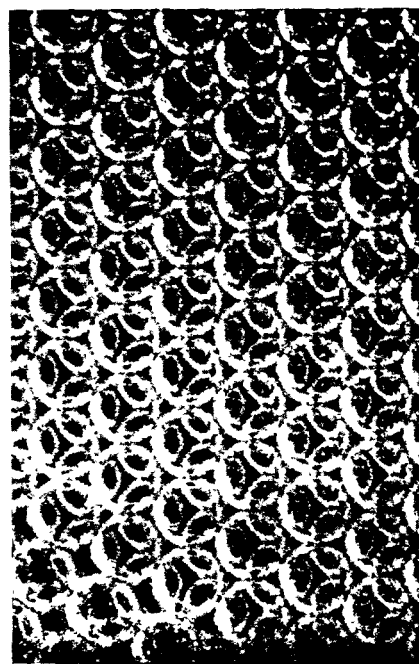


b. (100) face.

Face-centred cubic structure



c. Twin across (111), cubic structure.



d. Possible example of hexagonal close-packing

Diameter 0.70 mm.

FIGURE 17

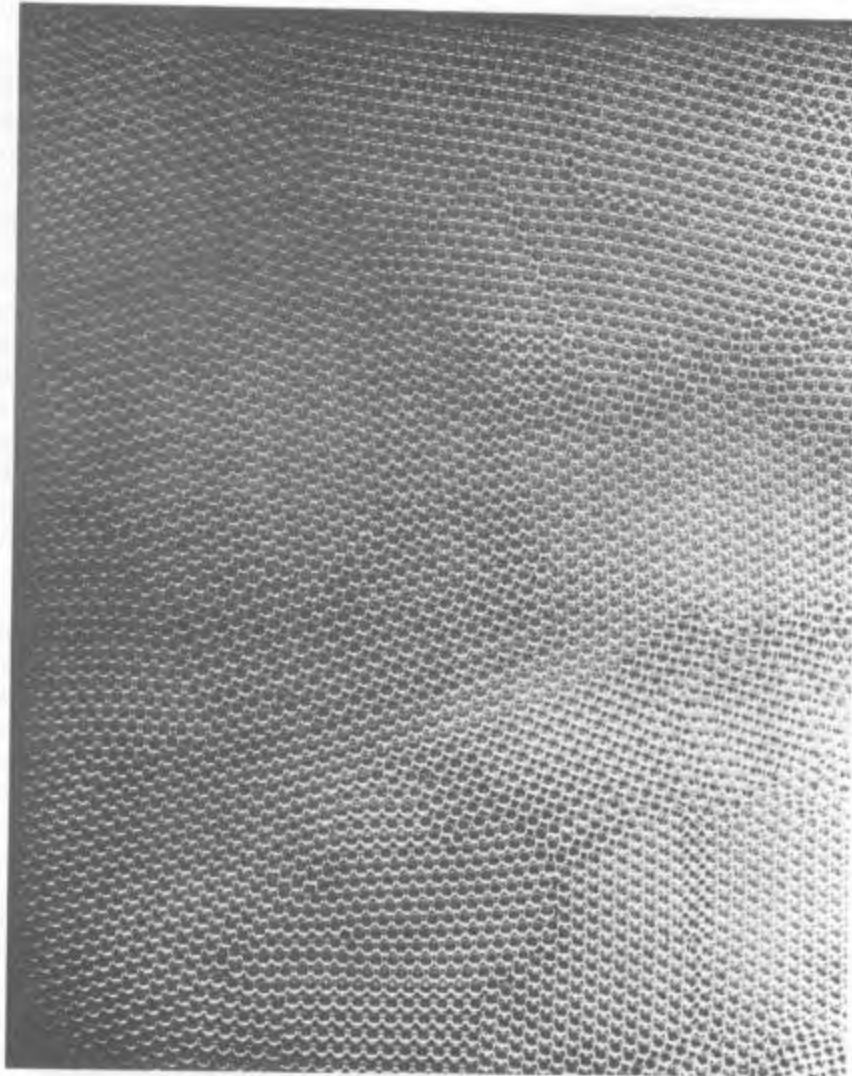


FIGURE 16. A three-dimensional raft viewed normally. Diameter 0.70 mm.

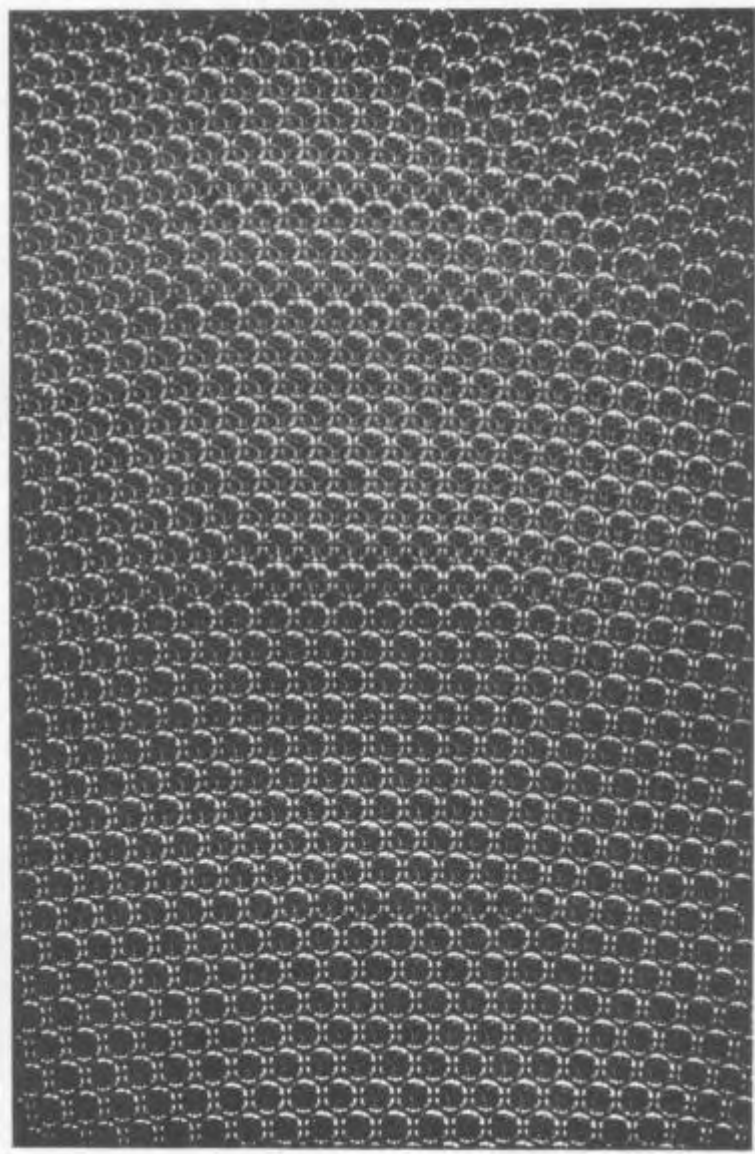


FIGURE 18. Dislocations in three-dimensional structure. Diameter 0.70 mm.

Tensors

31-1 The tensor of polarizability

Physicists always have a habit of taking the simplest example of any phenomenon and calling it “physics,” leaving the more complicated examples to become the concern of other fields—say of applied mathematics, electrical engineering, chemistry, or crystallography. Even solid-state physics is almost only half physics because it worries too much about special substances. So in these lectures we will be leaving out many interesting things. For instance, one of the important properties of crystals—or of most substances—is that their electric polarizability is different in different directions. If you apply a field in any direction, the atomic charges shift a little and produce a dipole moment, but the magnitude of the moment depends very much on the direction of the field. That is, of course, quite a complication. But in physics we usually start out by talking about the special case in which the polarizability is the same in all directions, to make life easier. We leave the other cases to some other field. Therefore, for our later work, we will not need at all what we are going to talk about in this chapter.

The mathematics of tensors is particularly useful for describing properties of substances which vary in direction—although that’s only one example of their use. Since most of you are not going to become physicists, but are going to go into the *real* world, where things depend severely upon direction, sooner or later you will need to use tensors. In order not to leave anything out, we are going to describe tensors, although not in great detail. We want the feeling that our treatment of physics is complete. For example, our electrodynamics is complete—as complete as any electricity and magnetism course, even a graduate course. Our mechanics is not complete, because we studied mechanics when you didn’t have a high level of mathematical sophistication, and we were not able to discuss subjects like the principle of least action, or Lagrangians, or Hamiltonians, and so on, which are *more elegant ways* of describing mechanics. Except for general relativity, however, we do have the complete *laws* of mechanics. Our electricity and magnetism is complete, and a lot of other things are quite complete. The quantum mechanics, naturally, will not be—we have to leave something for the future. But you should at least know what a tensor is.

We emphasized in Chapter 30 that the properties of crystalline substances are different in different directions—we say they are *anisotropic*. The variation of the induced dipole moment with the direction of the applied electric field is only one example, the one we will use for our example of a tensor. Let’s say that for a given direction of the electric field the induced dipole moment per unit volume \mathbf{P} is proportional to the strength of the applied field \mathbf{E} . (This is a good approximation for many substances if \mathbf{E} is not too large.) We will call the proportionality constant α .* We want now to consider substances in which α depends on the direction of the applied field, as, for example, in crystals like calcite, which make double images when you look through them.

Suppose, in a particular crystal, we find that an electric field \mathbf{E}_1 in the x -direction produces the polarization \mathbf{P}_1 in the x -direction. Then we find that an electric field \mathbf{E}_2 in the y -direction, with the same *strength*, as \mathbf{E}_1 produces a different polar-

31-1 The tensor of polarizability

31-2 Transforming the tensor components

31-3 The energy ellipsoid

31-4 Other tensors; the tensor of inertia

31-5 The cross product

31-6 The tensor of stress

31-7 Tensors of higher rank

31-8 The four-tensor of electromagnetic momentum

Review Chapter 11, Vol. I, *Vectors*
Chapter 20, Vol. I, *Rotation in Space*

* In Chapter 10 we followed the usual convention and wrote $\mathbf{P} = \epsilon_0 \chi \mathbf{E}$ and called χ (“khi”) the “susceptibility.” Here, it will be more convenient to use a single letter, so we write α for $\epsilon_0 \chi$. For isotropic dielectrics, $\alpha = (\kappa - 1)\epsilon_0$, where κ is the dielectric constant (see Section 10-4).

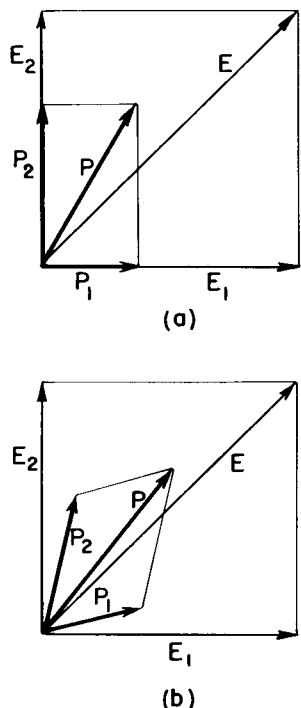


Fig. 31-1. The vector addition of polarizations in an anisotropic crystal.

ization P_2 in the y -direction. What would happen if we put an electric field at 45° ? Well, that's a superposition of two fields along x and y , so the polarization P will be the vector sum of P_1 and P_2 , as shown in Fig. 31-1(a). The polarization is no longer in the same direction as the electric field. You can see how that might come about. There may be charges which can move easily up and down, but which are rather stiff for sidewise motions. When a force is applied at 45° , the charges move farther up than they do toward the side. The displacements are not in the direction of the external force, because there are asymmetric internal elastic forces.

There is, of course, nothing special about 45° . It is *generally* true that the induced polarization of a crystal is *not* in the direction of the electric field. In our example above, we happened to make a "lucky" choice of our x - and y -axes, for which P was along E for both the x - and y -directions. If the crystal were rotated with respect to the coordinate axes, the electric field E_2 in the y -direction would have produced a polarization P with both an x - and a y -component. Similarly, the polarization due to an electric field in the x -direction would have produced a polarization with an x -component and a y -component. Then the polarizations would be as shown in Fig. 31-1(b), instead of as in part (a). Things get more complicated—but for any field E , the *magnitude* of P is still proportional to the magnitude of E .

We want now to treat the general case of an arbitrary orientation of a crystal with respect to the coordinate axes. An electric field in the x -direction will produce a polarization P with x -, y -, and z -components; we can write

$$P_x = \alpha_{xx}E_x, \quad P_y = \alpha_{yx}E_x, \quad P_z = \alpha_{zx}E_x. \quad (31.1)$$

All we are saying here is that if the electric field is in the x -direction, the polarization does not have to be in that same direction, but rather has an x -, a y -, and a z -component—each proportional to E_x . We are calling the constants of proportionality α_{xx} , α_{yx} , and α_{zx} , respectively (the first letter to tell us which component of P is involved, the last to refer to the direction of the electric field).

Similarly, for a field in the y -direction, we can write

$$P_x = \alpha_{xy}E_y, \quad P_y = \alpha_{yy}E_y, \quad P_z = \alpha_{zy}E_y; \quad (31.2)$$

and for a field in the z -direction,

$$P_x = \alpha_{xz}E_z, \quad P_y = \alpha_{yz}E_z, \quad P_z = \alpha_{zz}E_z. \quad (31.3)$$

Now we have said that polarization depends linearly on the fields, so if there is an electric field E that has both an x - and a y -component, the resulting x -component of P will be the sum of the two P_x 's of Eqs. (31.1) and (31.2). If E has components along x , y , and z , the resulting components of P will be the sum of the three contributions in Eqs. (31.1), (31.2), and (31.3). In other words, P will be given by

$$\begin{aligned} P_x &= \alpha_{xx}E_x + \alpha_{xy}E_y + \alpha_{xz}E_z, \\ P_y &= \alpha_{yx}E_x + \alpha_{yy}E_y + \alpha_{yz}E_z, \\ P_z &= \alpha_{zx}E_x + \alpha_{zy}E_y + \alpha_{zz}E_z. \end{aligned} \quad (31.4)$$

The dielectric behavior of the crystal is then completely described by the nine quantities (α_{xx} , α_{xy} , α_{xz} , α_{yx} , \dots), which we can represent by the symbol α_{ij} . (The subscripts i and j each stand for any one of the three possible letters x , y , and z .) Any arbitrary electric field E can be resolved with the components E_x , E_y , and E_z ; from these we can use the α_{ij} to find P_x , P_y , and P_z , which together give the total polarization P . The set of nine coefficients α_{ij} is called a *tensor*—in this instance, the *tensor of polarizability*. Just as we say that the three numbers (E_x , E_y , E_z) "form the vector E ," we say that the nine numbers (α_{xx} , α_{xy} , \dots) "form the tensor α_{ij} ."

31-2 Transforming the tensor components

You know that when we change to a different coordinate system x' , y' , and z' , the components $E_{x'}$, $E_{y'}$, and $E_{z'}$ of the vector will be quite different—as will also the components of \mathbf{P} . So all the coefficients α_{ij} will be different for a different set of coordinates. You can, in fact, see how the α 's must be changed by changing the components of \mathbf{E} and \mathbf{P} in the proper way, because if we describe the *same physical* electric field in the new coordinate system we should get the same polarization. For any new set of coordinates, $P_{x'}$ is a linear combination of P_x , P_y , and P_z :

$$P_{x'} = aP_x + bP_y + cP_z,$$

and similarly for the other components. If you substitute for P_x , P_y , and P_z in terms of the E 's, using Eq. (31.4), you get

$$\begin{aligned} P_{x'} &= a(\alpha_{xx}E_x + \alpha_{xy}E_y + \alpha_{xz}E_z) \\ &\quad + b(\alpha_{yx}E_x + \alpha_{yy}E_y + \cdots) \\ &\quad + c(\alpha_{zx}E_x + \cdots + \cdots). \end{aligned}$$

Then you write E_x , E_y , and E_z in terms of $E_{x'}$, $E_{y'}$, and $E_{z'}$; for instance,

$$E_x = a'E_{x'} + b'E_{y'} + c'E_{z'},$$

where a' , b' , c' are related to, but not equal to, a , b , c . So you have $P_{x'}$, expressed in terms of the components $E_{x'}$, $E_{y'}$, and $E_{z'}$; that is, you have the new α_{ij} . It is fairly messy, but quite straightforward.

When we talk about changing the axes we are assuming that the crystal stays put *in space*. If the crystal were rotated *with* the axes, the α 's would not change. Conversely, if the orientation of the crystal were changed with respect to the axes, we would have a new set of α 's. But if they are known for *any* one orientation of the crystal, they can be found for any other orientation by the transformation we have just described. In other words, the dielectric property of a crystal is described *completely* by giving the components of the polarization tensor α_{ij} , with respect to any arbitrarily chosen set of axes. Just as we can associate a vector velocity $\mathbf{v} = (v_x, v_y, v_z)$ with a particle, knowing that the three components will change in a certain definite way if we change our coordinate axes, so with a crystal we associate its polarization tensor α_{ij} , whose nine components will transform in a certain definite way if the coordinate system is changed.

The relation between \mathbf{P} and \mathbf{E} written in Eq. (31.4) can be put in the more compact notation:

$$P_i = \sum_j \alpha_{ij} E_j, \quad (31.5)$$

where it is understood that i represents either x , y , or z and that the sum is taken on $j = x, y$, and z . Many special notations have been invented for dealing with tensors, but each of them is convenient only for a limited class of problems. One common convention is to omit the sum sign (\sum) in Eq. (31.5), leaving it *understood* that whenever the same subscript occurs twice (here j), a sum is to be taken over that index. Since we will be using tensors so little, we will not bother to adopt any such special notations or conventions.

31-3 The energy ellipsoid

We want now to get some experience with tensors. Suppose we ask the interesting question: What energy is required to polarize the crystal (in addition to the energy in the electric field which we know is $\epsilon_0 E^2/2$ per unit volume)? Consider for a moment the atomic charges that are being displaced. The work done in displacing the charge the distance dx is $qE_x dx$, and if there are N charges per unit volume, the work done is $qE_x N dx$. But $qN dx$ is the change dP_x in the dipole

moment per unit volume. So the energy required *per unit volume* is

$$E_x dP_x.$$

Combining the work for the three components of the field, the work per unit volume is found to be

$$\mathbf{E} \cdot d\mathbf{P}.$$

Since the magnitude of \mathbf{P} is proportional to \mathbf{E} , the work done per unit volume in bringing the polarization from 0 to \mathbf{P} is the integral of $\mathbf{E} \cdot d\mathbf{P}$. Calling this work u_P ,* we write

$$u_P = \frac{1}{2} \mathbf{E} \cdot \mathbf{P} = \frac{1}{2} \sum_i E_i P_i. \quad (31.6)$$

Now we can express \mathbf{P} in terms of \mathbf{E} by Eq. (31.5), and we have that

$$u_P = \frac{1}{2} \sum_i \sum_j \alpha_{ij} E_i E_j. \quad (31.7)$$

The energy density u_P is a number independent of the choice of axes, so it is a scalar. A tensor has then the property that when it is summed over one index (with a vector), it gives a new vector; and when it is summed over *both* indexes (with *two* vectors), it gives a scalar.

The tensor α_{ij} should really be called a “tensor of second rank,” because it has two indexes. A vector—with *one* index—is a tensor of the first rank, and a scalar—with no index—is a tensor of zero rank. So we say that the electric field \mathbf{E} is a tensor of the first rank and that the energy density u_P is a tensor of zero rank. It is possible to extend the ideas of a tensor to three or more indexes, and so to make tensors of ranks higher than two.

The subscripts of the polarization tensor range over three possible values—they are tensors in three dimensions. The mathematicians consider tensors in four, five, or more dimensions. We have already used a four-dimensional tensor $F_{\mu\nu}$ in our relativistic description of the electromagnetic field (Chapter 26).

The polarization tensor α_{ij} has the interesting property that it is *symmetric*, that is, that $\alpha_{xy} = \alpha_{yx}$, and so on for any pair of indexes. (This is a *physical* property of a real crystal and not necessary for all tensors.) You can prove for yourself that this must be true by computing the change in energy of a crystal through the following cycle: (1) Turn on a field in the x -direction; (2) turn on a field in the y -direction; (3) turn *off* the x -field; (4) turn off the y -field. The crystal is now back where it started, and the net work done on the polarization must be back to zero. You can show, however, that for this to be true, α_{xy} must be equal to α_{yx} . The same kind of argument can, of course, be given for α_{xz} , etc. So the polarization tensor is symmetric.

This also means that the polarization tensor can be measured by just measuring the energy required to polarize the crystal in various directions. Suppose we apply an \mathbf{E} -field with only an x - and a y -component; then according to Eq. (31.7),

$$u_P = \frac{1}{2} [\alpha_{xx} E_x^2 + (\alpha_{xy} + \alpha_{yx}) E_x E_y + \alpha_{yy} E_y^2]. \quad (31.8)$$

With an E_x alone, we can determine α_{xx} ; with an E_y alone, we can determine α_{yy} ; with both E_x and E_y , we get an extra energy due to the term with $(\alpha_{xy} + \alpha_{yx})$. Since the α_{xy} and α_{yx} are equal, this term is $2\alpha_{xy}$ and can be related to the energy.

The energy expression, Eq. (31.8), has a nice geometric interpretation. Suppose we ask what fields E_x and E_y correspond to some *given* energy density—say u_0 . That is just the mathematical problem of solving the equation

$$\alpha_{xx} E_x^2 + 2\alpha_{xy} E_x E_y + \alpha_{yy} E_y^2 = 2u_0.$$

This is a quadratic equation, so if we plot E_x and E_y , the solutions of this equation

* This work done in *producing* the polarization by an electric field is not to be confused with the potential energy $-\mathbf{p}_0 \cdot \mathbf{E}$ of a permanent dipole moment \mathbf{p}_0 .

are all the points on an ellipse (Fig. 31-2). (It must be an ellipse, rather than a parabola or a hyperbola, because the energy for any field is always positive and finite.) The vector E with components E_x and E_y can be drawn from the origin to the ellipse. So such an “energy ellipse” is a nice way of “visualizing” the polarization tensor.

If we now generalize to include all three components, the electric vector E in any direction required to give a unit energy density gives a point which will be on the surface of an ellipsoid, as shown in Fig. 31-3. The shape of this ellipsoid of constant energy uniquely characterizes the tensor polarizability.

Now an ellipsoid has the nice property that it can always be described simply by giving the directions of three “principal axes” and the diameters of the ellipse along these axes. The “principal axes” are the directions of the longest and shortest diameters and the direction at right angles to both. They are indicated by the axes a , b , and c in Fig. 31-3. With respect to these axes, the ellipsoid has the particularly simple equation

$$\alpha_{aa}E_a^2 + \alpha_{bb}E_b^2 + \alpha_{cc}E_c^2 = 2u_0.$$

So with respect to these axes, the dielectric tensor has only three components that are not zero: α_{aa} , α_{bb} , and α_{cc} . That is to say, no matter how complicated a crystal is, it is always possible to choose a set of axes (not necessarily the crystal axes) for which the polarization tensor has only three components. With such a set of axes, Eq. (31.4) becomes simply

$$P_a = \alpha_{aa}E_a, \quad P_b = \alpha_{bb}E_b, \quad P_c = \alpha_{cc}E_c. \quad (31.9)$$

An electric field along any one of the principal axes produces a polarization along the same axis, but the coefficients for the three axes may, of course, be different.

Often, a tensor is described by listing the nine coefficients in a table inside of a pair of brackets:

$$\begin{bmatrix} \alpha_{xx} & \alpha_{xy} & \alpha_{xz} \\ \alpha_{yx} & \alpha_{yy} & \alpha_{yz} \\ \alpha_{zx} & \alpha_{zy} & \alpha_{zz} \end{bmatrix}. \quad (31.10)$$

For the principal axes a , b , and c , only the diagonal terms are not zero; we say then that “the tensor is diagonal.” The complete tensor is

$$\begin{bmatrix} \alpha_{aa} & 0 & 0 \\ 0 & \alpha_{bb} & 0 \\ 0 & 0 & \alpha_{cc} \end{bmatrix}. \quad (31.11)$$

The important point is that any polarization tensor (in fact, any symmetric tensor of rank two in any number of dimensions) can be put in this form by choosing a suitable set of coordinate axes.

If the three elements of the polarization tensor in diagonal form are all equal, that is, if

$$\alpha_{aa} = \alpha_{bb} = \alpha_{cc} = \alpha, \quad (31.12)$$

the energy ellipsoid becomes a sphere, and the polarizability is the same in all directions. The material is isotropic. In the tensor notation,

$$\alpha_{ij} = \alpha \delta_{ij}, \quad (31.13)$$

where δ_{ij} is the *unit tensor*

$$\delta_{ij} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}. \quad (31.14)$$

That means, of course,

$$\begin{aligned} \delta_{ij} &= 1, & \text{if } i &= j; \\ \delta_{ij} &= 0, & \text{if } i &\neq j. \end{aligned} \quad (31.15)$$

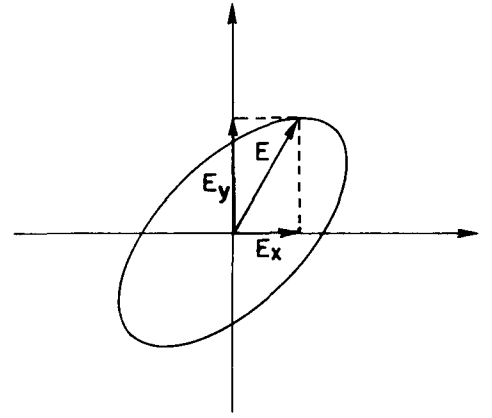


Fig. 31-2. Locus of the vector $E = (E_x, E_y)$ that gives a constant energy of polarization.

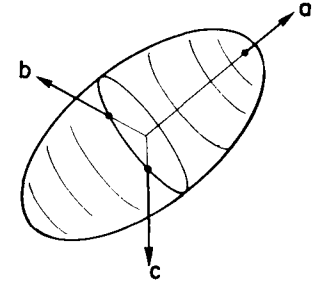


Fig. 31-3. The energy ellipsoid of the polarization tensor.

The tensor δ_{ij} is often called the “Kronecker delta.” You may amuse yourself by proving that the tensor (31.14) has exactly the same form if you change the coordinate system to any other rectangular one. The polarization tensor of Eq. (31.13) gives

$$P_i = \alpha \sum_j \delta_{ij} E_j = \alpha E_i,$$

which means the same as our old result for isotropic dielectrics:

$$\mathbf{P} = \alpha \mathbf{E}.$$

The shape and orientation of the polarization ellipsoid can sometimes be related to the symmetry properties of the crystal. We have said in Chapter 30 that there are 230 different possible internal symmetries of a three-dimensional lattice and that they can, for many purposes, be conveniently grouped into seven classes, according to the shape of the unit cell. Now the ellipsoid of polarizability must share the internal geometric symmetries of the crystal. For example, a triclinic crystal has low symmetry—the ellipsoid of polarizability will have unequal axes, and its orientation will not, in general, be aligned with the crystal axes. On the other hand, a monoclinic crystal has the property that its properties are unchanged if the crystal is rotated 180° about one axis. So the polarization tensor must be the same after such a rotation. It follows that the ellipsoid of the polarizability must return to itself after a 180° rotation. That can happen only if one of the axes of the ellipsoid is in the same direction as the symmetry axis of the crystal. Otherwise, the orientation and dimensions of the ellipsoid are unrestricted.

For an orthorhombic crystal, however, the axes of the ellipsoid must correspond to the crystal axes, because a 180° rotation about any one of the three axes repeats the same lattice. If we go to a tetragonal crystal, the ellipse must have the same symmetry, so it must have two equal diameters. Finally, for a cubic crystal, all three diameters of the ellipsoid must be equal, it becomes a sphere, and the polarizability of the crystal is the same in all directions.

There is a big game of figuring out the possible kinds of tensors for all the possible symmetries of a crystal. It is called a “group-theoretical” analysis. But for the simple case of the polarizability tensor, it is relatively easy to see what the relations must be.

31-4 Other tensors; the tensor of inertia

There are many other examples of tensors appearing in physics. For example, in a metal, or in any conductor, one often finds that the current density \mathbf{j} is approximately proportional to the electric field \mathbf{E} ; the proportionality constant is called the conductivity σ :

$$\mathbf{j} = \sigma \mathbf{E}.$$

For crystals, however, the relation between \mathbf{j} and \mathbf{E} is more complicated; the conductivity is not the same in all directions. The conductivity is a tensor, and we write

$$j_i = \sum_j \sigma_{ij} E_j.$$

Another example of a physical tensor is the moment of inertia. In Chapter 18 of Volume I we saw that a solid object rotating about a fixed axis has an angular momentum \mathbf{L} proportional to the angular velocity $\boldsymbol{\omega}$, and we called the proportionality factor I , the moment of inertia:

$$\mathbf{L} = I \boldsymbol{\omega}.$$

For an arbitrarily shaped object, the moment of inertia depends on its orientation with respect to the axis of rotation. For instance, a rectangular block will have different moments about each of its three orthogonal axes. Now angular velocity $\boldsymbol{\omega}$ and angular momentum \mathbf{L} are both vectors. For rotations about one of the axes of symmetry, they are parallel. But if the moment of inertia is different for the

three principal axes, then ω and L are, in general, not in the same direction (see Fig. 31-4). They are related in a way analogous to the relation between E and P . In general, we must write

$$\begin{aligned} L_x &= I_{xx}\omega_x + I_{xy}\omega_y + I_{xz}\omega_z, \\ L_y &= I_{xy}\omega_x + I_{yy}\omega_y + I_{yz}\omega_z, \\ L_z &= I_{zx}\omega_x + I_{zy}\omega_y + I_{zz}\omega_z. \end{aligned} \quad (31.16)$$

The nine coefficients I_{ij} are called the tensor of inertia. Following the analogy with the polarization, the kinetic energy for any angular momentum must be some quadratic form in the components ω_x , ω_y , and ω_z :

$$\text{KE} = \frac{1}{2} \sum_{ij} I_{ij} \omega_i \omega_j. \quad (31.17)$$

We can use the energy to define the ellipsoid of inertia. Also, energy arguments can be used to show that the tensor is symmetric—that $I_{ij} = I_{ji}$.

The tensor of inertia for a rigid body can be worked out if the shape of the object is known. We need only to write down the total kinetic energy of all the particles in the body. A particle of mass m and velocity v has the kinetic energy $\frac{1}{2}mv^2$, and the total kinetic energy is just the sum

$$\sum \frac{1}{2}mv^2$$

over all of the particles of the body. The velocity v of each particle is related to the angular velocity ω of the solid body. Let's assume that the body is rotating about its center of mass, which we take to be at rest. Then if r is the displacement of a particle from the center of mass, its velocity v is given by $\omega \times r$. So the total kinetic energy is

$$\text{KE} = \sum \frac{1}{2}m(\omega \times r)^2. \quad (31.18)$$

Now all we have to do is write $\omega \times r$ out in terms of the components ω_x , ω_y , ω_z , and x , y , z , and compare the result with Eq. (31.17); we find I_{ij} by identifying terms. Carrying out the algebra, we write

$$\begin{aligned} (\omega \times r)^2 &= (\omega \times r)_x^2 + (\omega \times r)_y^2 + (\omega \times r)_z^2 \\ &= (\omega_y z - \omega_z y)^2 + (\omega_z x - \omega_x z)^2 + (\omega_x y - \omega_y x)^2 \\ &= +\omega_y^2 z^2 - 2\omega_y \omega_z zy + \omega_z^2 y^2 \\ &\quad + \omega_z^2 x^2 - 2\omega_z \omega_x xz + \omega_x^2 z^2 \\ &\quad + \omega_x^2 y^2 - 2\omega_x \omega_y yx + \omega_y^2 x^2. \end{aligned}$$

Multiplying this equation by $m/2$, summing over all particles, and comparing with Eq. (31.17), we see that I_{xx} , for instance, is given by

$$I_{xx} = \sum m(y^2 + z^2).$$

This is the formula we have had before (Chapter 19, Vol. I) for the moment of inertia of a body about the x -axis. Since $r^2 = x^2 + y^2 + z^2$, we can also write this term as

$$I_{xx} = \sum m(r^2 - x^2).$$

Working out all of the other terms, the tensor of inertia can be written as

$$I_{ij} = \begin{bmatrix} \sum m(r^2 - x^2) & -\sum mxy & -\sum mxz \\ -\sum myx & \sum m(r^2 - y^2) & -\sum myz \\ -\sum mzx & -\sum mzy & \sum m(r^2 - z^2) \end{bmatrix}. \quad (31.19)$$

If you wish, this may be written in "tensor notation" as

$$I_{ij} = \sum m(r^2 \delta_{ij} - r_i r_j). \quad (31.20)$$

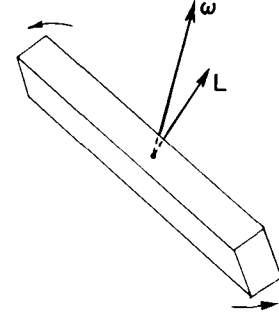


Fig. 31-4. The angular momentum L of a solid object is not, in general, parallel to its angular velocity ω .

where the r_i are the components (x, y, z) of the position vector of a particle and the \sum means to sum over all the particles. The moment of inertia, then, is a tensor of the second rank whose terms are a property of the body and relate \mathbf{L} to $\boldsymbol{\omega}$ by

$$L_i = \sum_j I_{ij} \omega_j \quad (31.21)$$

For a body of any shape whatever, we can find the ellipsoid of inertia and, therefore, the three principal axes. Referred to these axes, the tensor will be diagonal, so for any object there are always three orthogonal axes for which the angular velocity and angular momentum are parallel. They are called the principal axes of inertia.

31-5 The cross product

We should point out that we have been using tensors of the second rank since Chapter 20 of Volume I. There, we defined a “torque in a plane,” such as τ_{xy} , by

$$\tau_{xy} = xF_y - yF_x.$$

Generalized to three dimensions, we could write

$$\tau_{ij} = r_i F_j - r_j F_i. \quad (31.22)$$

The quantity τ_{ij} is a tensor of the second rank. One way to see that this is so is by combining τ_{ij} with some vector, say the unit vector \mathbf{e} , according to

$$\sum_j \tau_{ij} e_j.$$

If this quantity is a *vector*, then τ_{ij} must transform as a tensor—this is our definition of a tensor. Substituting for τ_{ij} , we have

$$\begin{aligned} \sum_j \tau_{ij} e_j &= \sum_j r_i F_j e_j - \sum_j r_j e_j F_i \\ &= r_i (\mathbf{F} \cdot \mathbf{e}) - (\mathbf{r} \cdot \mathbf{e}) F_i. \end{aligned}$$

Since the dot products are scalars, the two terms on the right-hand side are vectors, and likewise their difference. So τ_{ij} is a tensor.

But τ_{ij} is a special kind of tensor; it is *antisymmetric*, that is,

$$\tau_{ij} = -\tau_{ji},$$

so it has only three nonzero terms— τ_{xy} , τ_{yz} , and τ_{zx} . We were able to show in Chapter 20 of Volume I that these three terms, almost “by accident,” transform like the three components of a vector, so that we could *define*

$$\boldsymbol{\tau} = (\tau_x, \tau_y, \tau_z) = (\tau_{yz}, \tau_{zx}, \tau_{xy})$$

We say “by accident,” because it happens only in three dimensions. In four dimensions, for instance, an antisymmetric tensor of the second rank has *six* nonzero terms and certainly cannot be replaced by a vector with *four* components.

Just as the axial vector $\boldsymbol{\tau} = \mathbf{r} \times \mathbf{F}$ is a tensor, so also is every cross product of two polar vectors—all the same arguments apply. By luck, however, they are also representable by vectors (really pseudovectors), so our mathematics has been made easier for us.

Mathematically, if \mathbf{a} and \mathbf{b} are any two vectors, the nine quantities $a_i b_j$ form a tensor (although it may have no useful physical purpose). Thus, for the position vector \mathbf{r} , $r_i r_j$ is a tensor, and since δ_{ij} is also, we see that the right side of Eq. (31.20) is indeed a tensor. Likewise Eq. (31.22) is a tensor, since the two terms on the right-hand side are tensors.

31-6 The tensor of stress

The symmetric tensors we have described so far arose as coefficients in relating one vector to another. We would like to look now at a tensor which has a different physical significance—the tensor of *stress*. Suppose we have a solid object with various forces on it. We say that there are various “stresses” inside, by which we mean that there are internal forces between neighboring parts of the material. We have talked a little about such stresses in a two-dimensional case when we considered the surface tension in a stretched diaphragm in Section 12-3. We will now see that the internal forces in the material of a three-dimensional body can be described in terms of a tensor.

Consider a body of some elastic material—say a block of jello. If we make a cut through the block, the material on each side of the cut will, in general, get displaced by the internal forces. Before the cut was made, there must have been forces between the two parts of the block that kept the material in place; we can define the stresses in terms of these forces. Suppose we look at an imaginary plane perpendicular to the x -axis—like the plane σ in Fig. 31-5—and ask about the force across a small area $\Delta y \Delta z$ in this plane. The material on the left of the area exerts the force $\Delta \mathbf{F}_1$ on the material to the right, as shown in part (b) of the figure. There is, of course, the opposite reaction force $-\Delta \mathbf{F}_1$ exerted on the material to the left of the surface. If the area is small enough, we expect that $\Delta \mathbf{F}_1$ is proportional to the area $\Delta y \Delta z$.

You are already familiar with one kind of stress—the pressure in a static liquid. There the force is equal to the pressure times the area and is at right angles to the surface element. For solids—also for viscous liquids in motion—the force need not be normal to the surface; there are *shear* forces in addition to pressures (positive or negative). (By a “shear” force we mean the *tangential* components of the force across a surface.) All three components of the force must be taken into account. Notice also that if we make our cut on a plane with some other orientation, the forces will be different. A complete description of the internal stress requires a tensor.

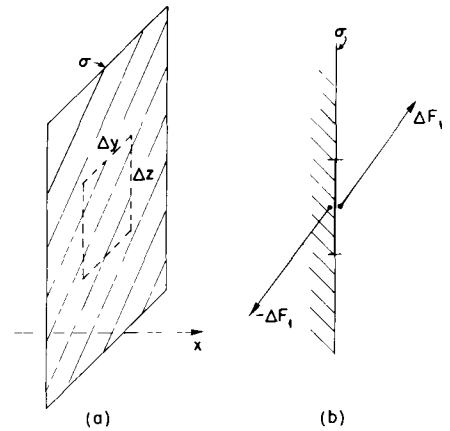


Fig. 31-5. The material to the left of the plane σ exerts across the area $\Delta y \Delta z$ the force $\Delta \mathbf{F}_1$ on the material to the right of the plane.

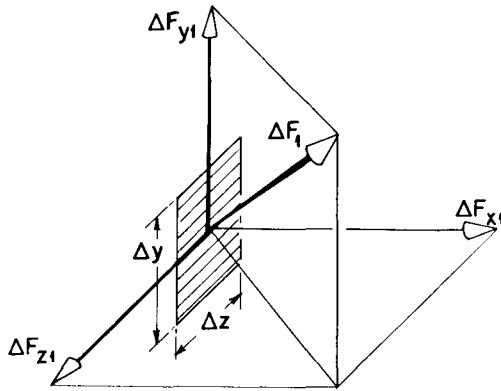


Fig. 31-6. The force $\Delta \mathbf{F}_1$ across an element of area $\Delta y \Delta z$ perpendicular to the x -axis is resolved into the three components ΔF_{x1} , ΔF_{y1} , and ΔF_{z1} .

We define the stress tensor in the following way: First, we imagine a cut perpendicular to the x -axis and resolve the force $\Delta \mathbf{F}_1$ across the cut into its components ΔF_{x1} , ΔF_{y1} , ΔF_{z1} , as in Fig. 31-6. The ratio of these forces to the area $\Delta y \Delta z$, we call S_{xx} , S_{yx} , and S_{zx} . For example,

$$S_{yx} = \frac{\Delta F_{y1}}{\Delta y \Delta z}.$$

The first index y refers to the direction force component; the second index x is normal to the area. If you wish, you can write the area $\Delta y \Delta z$ as Δa_x , meaning an element of area perpendicular to x . Then

$$S_{yx} = \frac{\Delta F_{y1}}{\Delta a_x}.$$

Next, we think of an imaginary cut perpendicular to the y -axis. Across a small

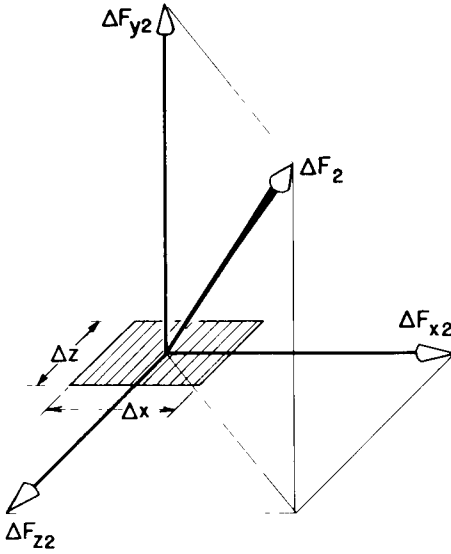


Fig. 31-7. The force across an element of area perpendicular to y is resolved into three rectangular components

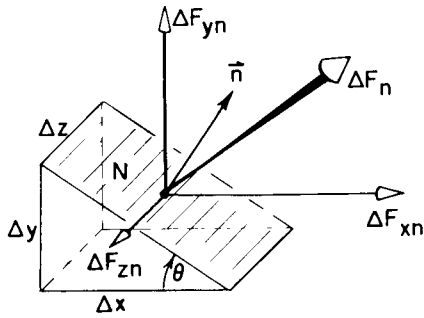


Fig. 31-8. The force F_n across the face N (whose unit normal is \mathbf{n}) is resolved into components.

area $\Delta x \Delta z$ there will be a force ΔF_2 . Again we resolve this force into three components, as shown in Fig. 31-7, and define the three components of the stress, S_{xy} , S_{yy} , S_{zy} , as the force per unit area in the three directions. Finally, we make an imaginary cut perpendicular to z and define the three components S_{xz} , S_{yz} , and S_{zz} . So we have the nine numbers

$$S_{ij} = \begin{bmatrix} S_{xx} & S_{xy} & S_{xz} \\ S_{yx} & S_{yy} & S_{yz} \\ S_{zx} & S_{zy} & S_{zz} \end{bmatrix}. \quad (31.23)$$

We want to show now that these nine numbers are sufficient to describe completely the internal state of stress, and that S_{ij} is indeed a tensor. Suppose we want to know the force across a surface oriented at some arbitrary angle. Can we find it from S_{ij} ? Yes, in the following way: We imagine a little solid figure which has one face N in the new surface, and the other faces parallel to the coordinate axes. If the face N happened to be parallel to the z -axis, we would have the triangular piece shown in Fig. 31-8. (This is a somewhat special case, but will illustrate well enough the general method.) Now the stress forces on the little solid triangle in Fig. 31-8 are in equilibrium (at least in the limit of infinitesimal dimensions), so the total force on it must be zero. We know the forces on the faces parallel to the coordinate axes directly from S_{ij} . Their vector sum must equal the force on the face N , so we can express this force in terms of S_{ij} .

Our assumption that the surface forces on the small triangular volume are in equilibrium neglects any other body forces that might be present, such as gravity or pseudo forces if our coordinate system is not an inertial frame. Notice, however, that such body forces will be proportional to the volume of the little triangle and, therefore, to Δx , Δy , Δz , whereas all the surface forces are proportional to the areas such as $\Delta x \Delta y$, $\Delta y \Delta z$, etc. So if we take the scale of the little wedge small enough, the body forces can always be neglected in comparison with the surface forces.

Let's now add up the forces on the little wedge. We take first the x -component, which is the sum of five parts—one from each face. However, if Δz is small enough, the forces on the triangular faces (perpendicular to the z -axis) will be equal and opposite, so we can forget them. The x -component of the force on the bottom rectangle is

$$\Delta F_{x2} = S_{xy} \Delta x \Delta z.$$

The x -component of the force on the vertical rectangle is

$$\Delta F_{x1} = S_{xx} \Delta y \Delta z.$$

These two must be equal to the x -component of the force outward across the face N . Let's call \mathbf{n} the unit vector normal to the face N , and the force on it \mathbf{F}_n , then we have

$$\Delta F_{xn} = S_{xx} \Delta y \Delta z + S_{xy} \Delta x \Delta z.$$

The x -component S_{xn} of the stress across this plane is equal to ΔF_{xn} divided by the area, which is $\Delta \sqrt{\Delta x^2 + \Delta y^2}$, or

$$S_{xn} = S_{xx} \frac{\Delta y}{\sqrt{\Delta x^2 + \Delta y^2}} + S_{xy} \frac{\Delta x}{\sqrt{\Delta x^2 + \Delta y^2}}.$$

Now $\Delta x / \sqrt{\Delta x^2 + \Delta y^2}$ is the cosine of the angle θ between \mathbf{n} and the y -axis, as shown in Fig. 31-8, so it can also be written as n_y , the y -component of \mathbf{n} . Similarly, $\Delta y / \sqrt{\Delta x^2 + \Delta y^2}$ is $\sin \theta = n_x$. We can write

$$S_{xn} = S_{xx} n_x + S_{xy} n_y.$$

If we now generalize to an arbitrary surface element, we would get that

$$S_{xn} = S_{xx} n_x + S_{xy} n_y + S_{xz} n_z$$

or, in general,

$$S_{in} = \sum_j S_{ij} n_j. \quad (31.24)$$

We can find the force across any surface element in terms of the S_{ij} , so it does describe completely the state of internal stress of the material.

Equation (31.24) says that the tensor S_{ij} relates the force S_n to the unit vector n , just as α_{ij} relates P to E . Since n and S_n are vectors, the components of S_{ij} must transform as a tensor with changes in coordinate axes. So S_{ij} is indeed a tensor.

We can also show that S_{ij} is a *symmetric* tensor by looking at the forces on a little cube of material. Suppose we take a little cube, oriented with its faces parallel to our coordinate axes, and look at it in cross section, as shown in Fig. 31-9. If we let the edge of the cube be one unit, the x - and y -components of the forces on the faces normal to the x - and y -axes might be as shown in the figure. If the cube is small, the stresses do not change appreciably from one side of the cube to the opposite side, so the force components are equal and opposite as shown. Now there must be no torque on the cube, or it would start spinning. The total torque about the center is $(S_{yx} - S_{xy})$ (times the unit edge of the cube), and since the total is zero, S_{yx} is equal to S_{xy} , and the stress tensor is symmetric.

Since S_{ij} is a symmetric tensor, it can be described by an ellipsoid which will have three principal axes. For surfaces normal to these axes, the stresses are particularly simple—they correspond to pushes or pulls perpendicular to the surfaces. There are no shear forces along these faces. For *any* stress, we can always choose our axes so that the shear components are zero. If the ellipsoid is a sphere, there are only normal forces in *any* direction. This corresponds to a hydrostatic pressure (positive or negative). So for a hydrostatic pressure, the tensor is diagonal and all three components are equal; they are, in fact, just equal to the pressure p . We can write

$$S_{ij} = p \delta_{ij}. \quad (31.25)$$

The stress tensor—and also its ellipsoid—will, in general, vary from point to point in a block of material; to describe the whole block we need to give the value of each component of S_{ij} as a function of position. So the stress tensor is a *field*. We have had *scalar fields*, like the temperature $T(x, y, z)$, which give one number for each point in space, and *vector fields* like $E(x, y, z)$, which give three numbers for each point. Now we have a *tensor field* which gives nine numbers for each point in space—or really six for the symmetric tensor S_{ij} . A complete description of the internal forces in an arbitrarily distorted solid requires six functions of x , y , and z .

31-7 Tensors of higher rank

The stress tensor S_{ij} describes the internal *forces* of matter. If the material is elastic, it is convenient to describe the internal *distortion* in terms of another tensor T_{ij} —called the *strain* tensor. For a simple object like a bar of metal, you know that the change in length, ΔL , is approximately proportional to the force, so we say it obeys Hooke's law:

$$\Delta L = \gamma F.$$

For a solid elastic body with arbitrary distortions, the strain T_{ij} is related to the stress S_{ij} by a set of linear equations:

$$T_{ij} = \sum_{k,l} \gamma_{ijkl} S_{kl}. \quad (31.26)$$

Also, you know that the potential energy of a spring (or bar) is

$$\frac{1}{2} F \Delta L = \frac{1}{2} \gamma F^2.$$

The generalization for the elastic energy *density* in a solid body is

$$U_{\text{elastic}} = \sum_{ijkl} \frac{1}{2} \gamma_{ijkl} S_{ij} S_{kl}. \quad (31.27)$$

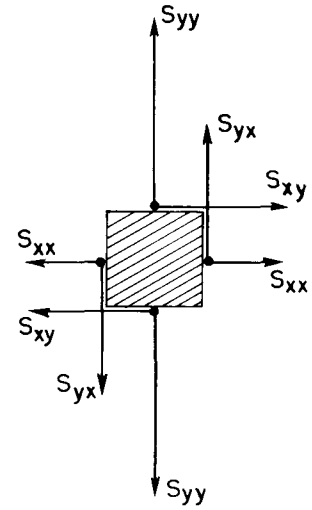


Fig. 31-9. The x - and y -forces on four faces of a small unit cube.

The complete description of the elastic properties of a crystal must be given in terms of the coefficients γ_{ijkl} . This introduces us to a new beast. It is a tensor of the *fourth* rank. Since each index can take on any one of three values, x, y , or z , there are $3^4 = 81$ coefficients. But there are really only 21 *different* numbers. First, since S_{ij} is symmetric, it has only six different values, and only 36 *different* coefficients are needed in Eq. (31.27). But also, S_{ij} can be interchanged with S_{kl} without changing the energy, so γ_{ijkl} must be symmetric if we interchange ij and kl . This reduces the number of different coefficients to 21. So to describe the elastic properties of a crystal of the lowest possible symmetry requires 21 elastic constants! This number is, of course, reduced for crystals of higher symmetry. For example, a cubic crystal has only three elastic constants, and an isotropic substance has only two.

That the latter is true can be seen as follows. How can the components of γ_{ijkl} be independent of the direction of the axes, as they must be if the material is isotropic? *Answer:* They can be independent *only* if they are expressible in terms of the tensor δ_{ij} . There are two possible expressions, $\delta_{ij}\delta_{kl}$ and $\delta_{ik}\delta_{jl} + \delta_{il}\delta_{jk}$, which have the required symmetry, so γ_{ijkl} must be a linear combination of them. Therefore, for isotropic materials,

$$\gamma_{ijkl} = a(\delta_{ij}\delta_{kl}) + b(\delta_{ik}\delta_{jl} + \delta_{il}\delta_{jk}),$$

and the material requires two constants, a and b , to describe its elastic properties. We will leave it for you to show that a cubic crystal needs only three

As a final example, this time of a third-rank tensor, we have the piezoelectric effect. Under stress, a crystal generates an electric field proportional to the stress; hence, in general, the law is

$$E_i = \sum_{j,k} P_{ijk} S_{jk},$$

where E_i is the electric field, and the P_{ijk} are the piezoelectric coefficients—or the piezoelectric tensor. Can you show that if the crystal has a center of inversion (invariant under $x, y, z \rightarrow -x, -y, -z$) the piezoelectric coefficients are all zero?

31-8 The four-tensor of electromagnetic momentum

All the tensors we have looked at so far in this chapter relate to the three dimensions of space; they are defined to have a certain transformation property under spatial rotations. In Chapter 26 we had occasion to use a tensor in the four dimensions of relativistic space-time—the electromagnetic field tensor $F_{\mu\nu}$. The components of such a four-tensor transform under a Lorentz transformation of the coordinates in a special way that we worked out. (Although we did not do it that way, we could have considered the Lorentz transformation as a “rotation” in a four-dimensional “space” called Minkowski space; then the analogy with what we are doing here would have been clearer.)

As our last example, we want to consider another tensor in the four dimensions (t, x, y, z) of relativity theory. When we wrote the stress tensor, we defined S_{ij} as a component of a force across a unit area. But a force is equal to the time rate of change of a momentum. Therefore, instead of saying “ S_{xy} is the x -component of the force across a unit area perpendicular to y ,” we could equally well say, “ S_{xy} is the rate of flow of the x -component of momentum through a unit area perpendicular to y .” In other words, each term of S_{ij} also represents the flow of the i -component of momentum through a unit area perpendicular to the j -direction. These are pure space components, but they are parts of a “larger” tensor $S_{\mu\nu}$ in four dimensions (μ and $\nu = t, x, y, z$) containing additional components like S_{tx} , S_{yt} , S_{tt} , etc. We will now try to find the physical meaning of these extra components.

We know that the space components represent flow of momentum. We can get a clue on how to extend this to the time dimension by studying another kind of “flow”—the flow of electric charge. For the *scalar* quantity, charge, the rate of flow (per unit area perpendicular to the flow) is a space *vector*—the current density

vector \mathbf{j} . We have seen that the time component of this flow vector is the density of the stuff that is flowing. For instance, \mathbf{j} can be combined with a time component, $j_t = \rho$, the charge density, to make the four-vector $j_\mu = (\rho, \mathbf{j})$; that is, the μ in j_μ takes on the values t, x, y, z to mean “density, rate of flow in the x -direction, rate of flow in y , rate of flow in z ” of the scalar charge.

Now by analogy with our statement about the time component of the flow of a scalar quantity, we might expect that with S_{xx} , S_{xy} , and S_{xz} , describing the flow of the x -component of momentum, there should be a time component S_{xt} which would be the density of whatever is flowing; that is, S_{xt} should be the density of x -momentum. So we can extend our tensor horizontally to include a t -component. We have

$$\begin{aligned} S_{xt} &= \text{density of } x\text{-momentum,} \\ S_{xx} &= x\text{-flow of } x\text{-momentum,} \\ S_{xy} &= y\text{-flow of } x\text{-momentum,} \\ S_{xz} &= z\text{-flow of } x\text{-momentum.} \end{aligned}$$

Similarly, for the y -component of momentum we have the three components of flow— S_{yx} , S_{yy} , S_{yz} —to which we should add a fourth term:

$$S_{yt} = \text{density of } y\text{-momentum.}$$

And, of course, to S_{zx} , S_{zy} , S_{zz} we would add

$$S_{zt} = \text{density of } z\text{-momentum.}$$

In four dimensions there is also a t -component of momentum, which is, we know, energy. So the tensor S_{ij} should be extended vertically with S_{tx} , S_{ty} , and S_{tz} , where

$$\begin{aligned} S_{tx} &= x\text{-flow of energy,} \\ S_{ty} &= y\text{-flow of energy,} \\ S_{tz} &= z\text{-flow of energy;} \end{aligned} \tag{31.28}$$

that is, S_{tx} is the flow of energy per unit area and per unit time across a surface perpendicular to the x -axis, and so on. Finally, to complete our tensor we need S_{tt} , which would be the *density of energy*. We have extended our stress tensor S_{ij} of three dimensions to the four-dimensional *stress-energy tensor* $S_{\mu\nu}$. The index μ can take on the four values t, x, y , and z , meaning, respectively, “density,” “flow per unit area in the x -direction,” “flow per unit area in the y -direction,” and “flow per unit area in the z -direction.” In the same way, ν takes on the four values t, x, y, z to tell us *what* flows, namely, “energy,” “momentum in the x -direction,” “momentum in the y -direction,” and “momentum in the z -direction.”

As an example, we will discuss this tensor not in matter, but in a region of free space in which there is an electromagnetic field. We know that the flow of energy is the Poynting vector $\mathbf{S} = \epsilon_0 c^2 \mathbf{E} \times \mathbf{B}$. So the x -, y -, and z -components of \mathbf{S} are, from the relativistic point of view, the components S_{tx} , S_{ty} , and S_{tz} of our four-dimensional stress-energy tensor. The symmetry of the tensor S_{ij} carries over into the time components as well, so the four-dimensional tensor $S_{\mu\nu}$ is symmetric:

$$S_{\mu\nu} = S_{\nu\mu}. \tag{31.29}$$

In other words, the components S_{xt} , S_{yt} , S_{zt} , which are the *densities* of x , y , and z momentum, are also equal to the x -, y -, and z -components of the Poynting vector \mathbf{S} , the *energy flow*—as we have already shown in an earlier chapter by a different kind of argument.

The remaining components of the electromagnetic stress tensor $S_{\mu\nu}$ can also be expressed in terms of the electric and magnetic fields \mathbf{E} and \mathbf{B} . That is to say, we must admit stress or, to put it less mysteriously, flow of momentum in the electromagnetic field. We discussed this in Chapter 27 in connection with Eq (27.21), but did not work out the details

Those who want to exercise their prowess in tensors in four dimensions might like to see the formula for $S_{\mu\nu}$ in terms of the fields:

$$S_{\mu\nu} = \frac{\epsilon_0}{2} \left(\sum_{\alpha} F_{\mu\alpha} F_{\nu\alpha} - \frac{1}{4} \delta_{\mu\nu} \sum_{\alpha,\beta} F_{\beta\alpha} F_{\beta\alpha} \right),$$

where sums on α, β are on t, x, y, z but (as usual in relativity) we adopt a special meaning for the sum sign \sum and for the symbol δ . In the sums the x, y, z terms are to be *subtracted* and $\delta_{tt} = +1$, while $\delta_{xx} = \delta_{yy} = \delta_{zz} = -1$ and $\delta_{\mu\nu} = 0$ for $\mu \neq \nu$ ($c = 1$). Can you verify that it gives the energy density $S_{tt} = (\epsilon_0/2) (E^2 + B^2)$ and the Poynting vector $\epsilon_0 \mathbf{E} \times \mathbf{B}$? Can you show that in an electrostatic field with $\mathbf{B} = 0$ the principal axes of stress are in the direction of the electric field, that there is a *tension* $(\epsilon_0/2)E^2$ along the direction of the field, and that there is an equal *pressure* in directions perpendicular to the field direction?

Refractive Index of Dense Materials

32-1 Polarization of matter

We want now to discuss the phenomenon of the refraction of light—and also, therefore, the absorption of light—by dense materials. In Chapter 31 of Volume I we discussed the theory of the index of refraction, but because of our limited mathematical abilities at that time, we had to restrict ourselves to finding the index only for materials of low density, like gases. The physical principles that produced the index were, however, made clear. The electric field of the light wave polarizes the molecules of the gas, producing oscillating dipole moments. The acceleration of the oscillating charges radiates new waves of the field. This new field, interfering with the old field, produces a changed field which is equivalent to a phase shift of the original wave. Because this phase shift is proportional to the thickness of the material, the effect is equivalent to having a different phase velocity in the material. When we looked at the subject before, we neglected the complications that arise from such effects as the new wave changing the fields at the oscillating dipoles. We assumed that the forces on the charges in the atoms came just from the *incoming* wave, whereas, in fact, their oscillations are driven not only by the incoming wave but also by the radiated waves of all the other atoms. It would have been difficult for us at that time to include this effect, so we studied only the rarefied gas, where such effects are not important.

Now, however, we will find that it is very easy to treat the problem by the use of differential equations. This method obscures the physical origin of the index (as coming from the re-radiated waves interfering with the original waves), but it makes the theory for dense materials much simpler. This chapter will bring together a large number of pieces from our earlier work. We've taken up practically everything we will need, so there are relatively few really new ideas to be introduced. Since you may need to refresh your memory about what we are going to need, we give in Table 32-1 a list of the equations we are going to use, together with a reference to the place where each can be found. In most instances, we will not take the time to give the physical arguments again, but will just use the equations.

Table 32-1

**Our work in this chapter will be based on the following material,
already covered in earlier chapters**

Subject	Reference	Equation
Damped oscillations	Vol. I, Chap. 23	$m(\ddot{x} + \gamma\dot{x} + \omega_0^2x) = F$
Index of gases	Vol. I, Chap. 31	$n = 1 + \frac{1}{2} \frac{Nq_e^2}{\epsilon_0(\omega_0^2 - \omega^2)}$ $n = n' - in''$
Mobility	Vol. I, Chap. 41	$m\dot{x} + \mu x = F$
Electrical conductivity	Vol. I, Chap. 43	$\mu = \frac{\tau}{m}; \sigma = \frac{Nq_e^2\tau}{m}$
Polarizability	Vol. II, Chap. 10	$\rho_{\text{pol}} = -\nabla \cdot \mathbf{P}$
Inside dielectrics	Vol. II, Chap. 11	$\mathbf{E}_{\text{local}} = \mathbf{E} + \frac{1}{3\epsilon_0} \mathbf{P}$

32-1 Polarization of matter

32-2 Maxwell's equations in a dielectric

32-3 Waves in a dielectric

32-4 The complex index of refraction

32-5 The index of a mixture

32-6 Waves in metals

32-7 Low-frequency and high-frequency approximations; the skin depth and the plasma frequency

Review: See Table 32-1.

We begin by recalling the machinery of the index of refraction for a gas. We suppose that there are N particles per unit volume and that each particle behaves as a harmonic oscillator. We use a model of an atom or molecule in which the electron is bound with a force proportional to its displacement (as though the electron were held in place by a spring). We emphasized that this was not a legitimate *classical* model of an atom, but we will show later that the correct quantum mechanical theory gives results equivalent to this model (in simple cases). In our earlier treatment, we did not include the possibility of a damping force in the atomic oscillators, but we will do so now. Such a force corresponds to a resistance to the motion, that is, to a force proportional to the velocity of the electron. Then the equation of motion is

$$F = q_e E = m(\ddot{x} + \gamma \dot{x} + \omega_0^2 x), \quad (32.1)$$

where x is the displacement parallel to the direction of E . (We are assuming an *isotropic* oscillator whose restoring force is the same in all directions. Also, we are taking, for the moment, a linearly polarized wave, so that E doesn't change direction.) If the electric field acting on the atom varies sinusoidally with time, we write

$$E = E_0 e^{i\omega t}. \quad (32.2)$$

The displacement will then oscillate with the same frequency, and we can let

$$x = x_0 e^{i\omega t}$$

Substituting $\dot{x} = i\omega x$ and $\ddot{x} = -\omega^2 x$, we can solve for x in terms of E :

$$x = \frac{q_e/m}{-\omega^2 + i\gamma\omega + \omega_0^2} E \quad (32.3)$$

Knowing the displacement, we can calculate the acceleration \ddot{x} and find the radiated wave responsible for the index. This was the way we computed the index in Chapter 31 of Volume I.

Now, however, we want to take a different approach. The induced dipole moment p of an atom is $q_e x$ or, using Eq. (32.3),

$$p = \frac{q_e^2/m}{-\omega^2 + i\gamma\omega + \omega_0^2} E. \quad (32.4)$$

Since p is proportional to E , we write

$$p = \epsilon_0 \alpha(\omega) E, \quad (32.5)$$

where α is called the *atomic polarizability*.* With this definition, we have

$$\alpha = \frac{q_e^2/m\epsilon_0}{-\omega^2 + i\gamma\omega + \omega_0^2}. \quad (32.6)$$

The quantum mechanical solution for the motions of electrons in atoms gives a similar answer except with the following modifications. The atoms have several natural frequencies, each frequency with its own dissipation constant γ . Also the effective "strength" of each mode is different, which we can represent by multiplying the polarizability for each frequency by a strength factor f , which is a number we expect to be of the order of 1. Representing the three parameters ω , γ , and f by ω_k , γ_k , and f_k for each mode of oscillation, and summing over the

* Throughout this chapter we follow the notation of Chapter 31 of Volume I, and let α represent the *atomic* polarizability as defined here. In the last chapter, we used α to represent the *volume* polarizability—the ratio of P to E . In the notation of *this* chapter $P = N\alpha\epsilon_0 E$ (see Eq. 32.8)

various modes, we modify Eq. (32.6) to read

$$\alpha(\omega) = \frac{q_e^2}{\epsilon_0 m} \sum_k \frac{f_k}{-\omega^2 + i\gamma_k \omega + \omega_{0k}^2}. \quad (32.7)$$

If N is the number of atoms per unit volume in the material, the polarization P is just $Np = \epsilon_0 N \alpha E$, and is proportional to E :

$$\mathbf{P} = \epsilon_0 N \alpha(\omega) \mathbf{E}. \quad (32.8)$$

In other words, when there is a sinusoidal electric field acting in a material, there is an induced dipole moment per unit volume which is proportional to the electric field—with a proportionality constant α that, we emphasize, depends upon the frequency. At very high frequencies, α is small; there is not much response. However, at low frequencies there can be a strong response. Also, the proportionality constant is a complex number, which means that the polarization does not exactly follow the electric field, but may be shifted in phase to some extent. At any rate, there is a polarization per unit volume whose magnitude is proportional to the strength of the electric field.

32-2 Maxwell's equations in a dielectric

The existence of polarization in matter means that there are polarization charges and currents inside of the material, and these must be put into the complete Maxwell equations in order to find the fields. We are going to solve Maxwell's equations this time in a situation in which the charges and currents are not zero, as in a vacuum, but are given implicitly by the polarization vector. Our first step is to find explicitly the charge density ρ and current density \mathbf{j} , averaged over a small volume of the same size we had in mind when we defined \mathbf{P} . Then the ρ and \mathbf{j} we need can be obtained from the polarization.

We have seen in Chapter 10 that when the polarization \mathbf{P} varies from place to place, there is a charge density given by

$$\rho_{\text{pol}} = -\nabla \cdot \mathbf{P}. \quad (32.9)$$

At that time, we were dealing with static fields, but the same formula is valid also for time-varying fields. However, when \mathbf{P} varies with time, there are charges in motion, so there is also a polarization *current*. Each of the oscillating charges contributes a current equal to its charge q_e , times its velocity v . With N such charges per unit volume, the current density \mathbf{j} is

$$\mathbf{j} = N q_e v.$$

Since we know that $v = dx/dt$, then $j = N q_e (dx/dt)$, which is just dP/dt . Therefore the current density from the varying polarization is

$$\mathbf{j}_{\text{pol}} = \frac{d\mathbf{P}}{dt} \quad (32.10)$$

Our problem is now direct and simple. We write Maxwell's equations with the charge density and current density expressed in terms of \mathbf{P} , using Eqs. (32.9) and (32.10). (We assume that there are no other currents and charges in the material.) We then relate \mathbf{P} to \mathbf{E} with Eq. (32.5), and we solve the equation for \mathbf{E} and \mathbf{B} —looking for the wave solutions.

Before we do this, we would like to make an historical note. Maxwell originally wrote his equations in a form which was different from the one we have been using. Because the equations were written in this different form for many years—and are still written that way by many people—we will explain the difference. In the early days, the mechanism of the dielectric constant was not fully and clearly appreciated. The nature of atoms was not understood, nor that there was a polarization of the material. So people did not appreciate that there was a contribution

to the charge density ρ from $\nabla \cdot \mathbf{P}$. They thought only in terms of charges that were not bound to atoms (such as the charges that flow in wires or are rubbed off surfaces).

Today, we prefer to let ρ represent the *total* charge density, including the part from the bound atomic charges. If we call that part ρ_{pol} , we can write

$$\rho = \rho_{\text{pol}} + \rho_{\text{other}},$$

where ρ_{other} is the charge density considered by Maxwell and refers to the charges not bound to individual atoms. We would then write

$$\nabla \cdot \mathbf{E} = \frac{\rho_{\text{pol}} + \rho_{\text{other}}}{\epsilon_0}.$$

Substituting ρ_{pol} from Eq. (32.9),

$$\nabla \cdot \mathbf{E} = \frac{\rho_{\text{other}}}{\epsilon_0} - \frac{1}{\epsilon_0} \nabla \cdot \mathbf{P}$$

or

$$\nabla \cdot (\epsilon_0 \mathbf{E} + \mathbf{P}) = \rho_{\text{other}}. \quad (32.11)$$

The current density in the Maxwell equations for $\nabla \times \mathbf{B}$ also has, in general, contributions from bound atomic currents. We can therefore write

$$\mathbf{j} = \mathbf{j}_{\text{pol}} + \mathbf{j}_{\text{other}},$$

and the Maxwell equation becomes

$$c^2 \nabla \times \mathbf{B} = \frac{\mathbf{j}_{\text{other}}}{\epsilon_0} + \frac{\mathbf{j}_{\text{pol}}}{\epsilon_0} + \frac{\partial \mathbf{E}}{\partial t}. \quad (32.12)$$

Using Eq. (32.10), we get

$$\epsilon_0 c^2 \nabla \times \mathbf{B} = \mathbf{j}_{\text{other}} + \frac{\partial}{\partial t} (\epsilon_0 \mathbf{E} + \mathbf{P}). \quad (32.13)$$

Now you can see that if we were to *define* a new vector \mathbf{D} by

$$\mathbf{D} = \epsilon_0 \mathbf{E} + \mathbf{P}, \quad (32.14)$$

the two field equations would become

$$\nabla \cdot \mathbf{D} = \rho_{\text{other}} \quad (32.15)$$

and

$$\epsilon_0 c^2 \nabla \times \mathbf{B} = \mathbf{j}_{\text{other}} + \frac{\partial \mathbf{D}}{\partial t}. \quad (32.16)$$

These are actually the forms that Maxwell used for dielectrics. His two remaining equations were

$$\nabla \times \mathbf{E} = - \frac{\partial \mathbf{B}}{\partial t},$$

and

$$\nabla \cdot \mathbf{B} = 0,$$

which are the same as we have been using.

Maxwell and the other early workers also had a problem with magnetic materials (which we will take up soon). Because they did not know about the circulating currents responsible for atomic magnetism, they used a current density that was missing still another part. Instead of Eq. (32.16), they actually wrote

$$\nabla \times \mathbf{H} = \mathbf{j}' + \frac{\partial \mathbf{D}}{\partial t}, \quad (32.17)$$

where \mathbf{H} differs from $\epsilon_0 c^2 \mathbf{B}$ because it includes the effects of atomic currents. (Then \mathbf{j}' represents what is left of the currents.) So Maxwell had *four* field vectors— \mathbf{E} , \mathbf{D} , \mathbf{B} , and \mathbf{H} —the \mathbf{D} and \mathbf{H} were hidden ways of not paying attention to what

was going on inside the material. You will find the equations written this way in many places.

To solve the equations, it is necessary to relate \mathbf{D} and \mathbf{H} to the other fields, and people used to write

$$\mathbf{D} = \epsilon \mathbf{E} \quad \text{and} \quad \mathbf{B} = \mu \mathbf{H}. \quad (32.18)$$

However, these relations are only approximately true for some materials and even then only if the fields are not changing rapidly with time. (For sinusoidally varying fields one often *can* write the equations this way by making ϵ and μ complex functions of the frequency, but not for an arbitrary time variation of the fields.) So there used to be all kinds of cheating in solving the equations. We think the right way is to keep the equations in terms of the fundamental quantities as we now understand them—and that's how we have done it.

32-3 Waves in a dielectric

We want now to find out what kind of electromagnetic waves can exist in a dielectric material in which there are no extra charges other than those bound in atoms. So we take $\rho = -\nabla \cdot \mathbf{P}$ and $\mathbf{j} = \partial \mathbf{P} / \partial t$. Maxwell's equations then become

$$\begin{aligned} \text{(a)} \quad \nabla \cdot \mathbf{E} &= -\frac{\nabla \cdot \mathbf{P}}{\epsilon_0} & \text{(b)} \quad c^2 \nabla \times \mathbf{B} &= \frac{\partial}{\partial t} \left(\frac{\mathbf{P}}{\epsilon_0} + \mathbf{E} \right) \\ \text{(c)} \quad \nabla \times \mathbf{E} &= -\frac{\partial \mathbf{B}}{\partial t} & \text{(d)} \quad \nabla \cdot \mathbf{B} &= 0 \end{aligned} \quad (32.19)$$

We can solve these equations as we have done before. We start by taking the curl of Eq. (32.19c):

$$\nabla \times (\nabla \times \mathbf{E}) = -\frac{\partial}{\partial t} \nabla \times \mathbf{B}.$$

Next, we make use of the vector identity

$$\nabla \times (\nabla \times \mathbf{E}) = \nabla(\nabla \cdot \mathbf{E}) - \nabla^2 \mathbf{E},$$

and also substitute for $\nabla \times \mathbf{B}$, using Eq. (32.19b); we get

$$\nabla(\nabla \cdot \mathbf{E}) - \nabla^2 \mathbf{E} = -\frac{1}{\epsilon_0 c^2} \frac{\partial^2 \mathbf{P}}{\partial t^2} - \frac{1}{c^2} \frac{\partial^2 \mathbf{E}}{\partial t^2}.$$

Using Eq. (32.19a) for $\nabla \cdot \mathbf{E}$, we get

$$\nabla^2 \mathbf{E} - \frac{1}{c^2} \frac{\partial^2 \mathbf{E}}{\partial t^2} = -\frac{1}{\epsilon_0 c^2} \nabla(\nabla \cdot \mathbf{P}) + \frac{1}{\epsilon_0 c^2} \frac{\partial^2 \mathbf{P}}{\partial t^2}. \quad (32.20)$$

So instead of the wave equation, we now get that the D'Alembertian of \mathbf{E} is equal to two terms involving the polarization \mathbf{P} .

Since \mathbf{P} depends on \mathbf{E} , however, Eq. (32.20) can still have wave solutions. We will now limit ourselves to *isotropic* dielectrics, so that \mathbf{P} is always in the same direction as \mathbf{E} . Let's try to find a solution for a wave going in the z -direction. Then, the electric field might vary as $e^{i(\omega t - kz)}$. We will also suppose that the wave is polarized in the x -direction—that the electric field has only an x -component. We write

$$E_x = E_0 e^{i(\omega t - kz)}. \quad (32.21)$$

You know that any function of $(z - vt)$ represents a wave that travels with the speed v . The exponent of Eq. (32.21) can be written as

$$-ik \left(z - \frac{\omega}{k} t \right),$$

so, Eq. (32.21) represents a wave with the phase velocity

$$v_{\text{ph}} = \omega/k.$$

The index of refraction n is defined (see Chapter 31, Vol. I) by letting

$$v_{\text{ph}} = \frac{c}{n}.$$

Thus Eq. (32.21) becomes

$$E_x = E_0 e^{i\omega(t - nz/c)}.$$

So we can find n by finding what value of k is required if Eq. (32.21) is to satisfy the proper field equations, and then using

$$n = \frac{kc}{\omega}. \quad (32.22)$$

In an isotropic material, there will be only an x -component of the polarization; then \mathbf{P} has no variation with the x -coordinate, so $\nabla \cdot \mathbf{P} = 0$, and we get rid of the first term on the right-hand side of Eq. (32.20). Also, since we are assuming a linear dielectric, P_x will vary as $e^{i\omega t}$, and $\partial^2 P_x / \partial t^2 = -\omega^2 P_x$. The Laplacian in Eq. (32.20) becomes simply $\partial^2 E_x / \partial z^2 = -k^2 E_x$, so we get

$$-k^2 E_x + \frac{\omega^2}{c^2} E_x = -\frac{\omega^2}{\epsilon_0 c^2} P_x. \quad (32.23)$$

Now let us assume for the moment that since \mathbf{E} is varying sinusoidally, we can set \mathbf{P} proportional to \mathbf{E} , as in Eq. (32.5). (We'll come back to discuss this assumption later.) We write

$$P_x = \epsilon_0 N \alpha E_x.$$

Then E_x drops out of Eq. (32.23), and we find

$$k^2 = \frac{\omega^2}{c^2} (1 + N\alpha). \quad (32.24)$$

We have found that a wave like Eq. (32.21), with the wave number k given by Eq. (32.24), will satisfy the field equations. Using Eq. (32.22), the index n is given by

$$n^2 = 1 + N\alpha. \quad (32.25)$$

Let's compare this formula with what we obtained in our theory of the index of a gas (Chapter 31, Vol. I). There, we got Eq. (31.29), which is

$$n = 1 + \frac{1}{2} \frac{Nq_e^2}{m\epsilon_0} \frac{1}{-\omega^2 + \omega_0^2}. \quad (32.26)$$

Taking α from Eq. (32.6), Eq. (32.25) would give us

$$n^2 = 1 + \frac{Nq_e^2}{m\epsilon_0} \frac{1}{-\omega^2 + i\gamma\omega + \omega_0^2}. \quad (32.27)$$

First, we have the new term in $i\gamma\omega$, because we are including the dissipation of the oscillators. Second, the left-hand side is n instead of n^2 , and there is an extra factor of $1/2$. But notice that if N is small enough so that n is close to one (as it is for a gas), then Eq. (32.27) says that n^2 is one plus a small number: $n^2 = 1 + \epsilon$. We can then write $n = \sqrt{1 + \epsilon} \approx 1 + \epsilon/2$, and the two expressions are equivalent. Thus our new method gives for a gas the same result we found earlier.

Now you might think that Eq. (32.27) should give the index of refraction for dense materials also. It needs to be modified, however, for several reasons. First, the derivation of this equation assumes that the polarizing field on each atom is the field E_x . That assumption is *not* right, however, because in dense materials there is also the field produced by other atoms in the vicinity, which may be comparable to E_x . We considered a similar problem when we studied the static fields in dielectrics. (See Chapter 11.) You will remember that we estimated the field at a single atom by imagining that it sat in a spherical hole in the surrounding dielectric. The field in such a hole—which we called the *local* field—is increased

over the average field E by the amount $P/3\epsilon_0$. (Remember, however, that this result is only strictly true in isotropic materials—including the special case of a cubic crystal.)

The same arguments will hold for the electric field in a wave, so long as the wavelength of the wave is much longer than the spacing between atoms. Limiting ourselves to such cases, we write

$$E_{\text{local}} = E + \frac{P}{3\epsilon_0}. \quad (32.28)$$

This local field is the one that should be used for E in Eq. (32.3); that is, Eq.(32.8) should be rewritten:

$$P = \epsilon_0 N \alpha E_{\text{local}}. \quad (32.29)$$

Using E_{local} from Eq. (32.28), we find

$$P = \epsilon_0 N \alpha \left(E + \frac{P}{3\epsilon_0} \right)$$

or

$$P = \frac{N\alpha}{1 - (N\alpha/3)} \epsilon_0 E. \quad (32.30)$$

In other words, for dense materials P is still proportional to E (for sinusoidal fields). However, the constant of proportionality is not $\epsilon_0 N \alpha$, as we wrote below Eq. (32.23), but should be $\epsilon_0 N \alpha / [1 - (N\alpha/3)]$. So we should correct Eq. (32.25) to read

$$n^2 = 1 + \frac{N\alpha}{1 - (N\alpha/3)}. \quad (32.31)$$

It will be more convenient if we rewrite this equation as

$$3 \frac{n^2 - 1}{n^2 + 2} = N\alpha, \quad (32.32)$$

which is algebraically equivalent. This is known as the Clausius-Mosotti equation.

There is another complication in dense materials. Because neighboring atoms are so close, there are strong interactions between them. The internal modes of oscillation are, therefore, modified. The natural frequencies of the atomic oscillations are spread out by the interactions, and they are usually quite heavily damped—the resistance coefficient becomes quite large. So the ω_0 's and γ 's of the solid will be quite different from those of the free atoms. With these reservations, we can still represent α , at least approximately, by Eq. (32.7). We have then that

$$3 \frac{n^2 - 1}{n^2 + 2} = \frac{Nq_e^2}{m\epsilon_0} \sum_k \frac{f_k}{-\omega^2 + i\gamma_k\omega + \omega_{0k}^2}. \quad (32.33)$$

One final complication. If the dense material is a mixture of several components, each will contribute to the polarization. The total α will be the sum of the contributions from each component of the mixture [except for the inaccuracy of the local field approximation, Eq. (32.28), in ordered crystals—effects we discussed when analyzing ferroelectrics]. Writing N_j as the number of atoms of each component per unit volume, we should replace Eq. (32.32) by

$$3 \left(\frac{n^2 - 1}{n^2 + 2} \right) = \sum_j N_j \alpha_j, \quad (32.34)$$

where each α_j will be given by an expression like Eq. (32.7). Equation (32.34) completes our theory of the index of refraction. The quantity $3(n^2 - 1)/(n^2 + 2)$ is given by some complex function of frequency, which is the mean atomic polarizability $\alpha(\omega)$. The precise evaluation of $\alpha(\omega)$ (that is, finding f_k , γ_k and ω_{0k}) in dense substances is a difficult problem of quantum mechanics. It has been done from first principles only for a few especially simple substances.

32-4 The complex index of refraction

We want to look now at the consequences of our result, Eq. (32.33). First, we notice that α is complex, so the index n is going to be a complex number. What does that mean? Let's say that we write n as the sum of a real and an imaginary part:

$$n = n_R - in_I, \quad (32.35)$$

where n_R and n_I are real functions of ω . We write in_I with a minus sign, so that n_I will be a positive quantity in all ordinary optical materials. (In ordinary inactive materials—that are not, like lasers, light sources themselves— γ is a positive number, and that makes the imaginary part of n negative.) Our plane wave of Eq. (32.21) is written in terms of n as

$$E_x = E_0 e^{-i\omega(t - nz/c)}$$

Writing n as in Eq. (32.35), we would have

$$E_x = E_0 e^{-\omega n_I z/c} e^{i\omega(t - n_R z/c)}. \quad (32.36)$$

The term $e^{i\omega(t - n_R z/c)}$ represents a wave travelling with the speed c/n_R , so n_R represents what we normally think of as the index of refraction. But the *amplitude* of this wave is

$$E_0 e^{-\omega n_I z/c},$$

which decreases exponentially with z . A graph of the strength of the electric field at some instant as a function of z is shown in Fig. 32-1, for $n_I \approx n_R/2\pi$. The imaginary part of the index represents the attenuation of the wave due to the energy losses in the atomic oscillators. The *intensity* of the wave is proportional to the square of the amplitude, so

$$\text{Intensity} \propto e^{-2\omega n_I z/c}.$$

This is often written as

$$\text{Intensity} \propto e^{-\beta z},$$

where $\beta = 2\omega n_I/c$ is called the *absorption coefficient*. Thus we have in Eq. (32.33) not only the theory of the index of refraction of materials, but the theory of their absorption of light as well.

In what we usually consider to be transparent material, the quantity $c/\omega n_I$ —which has the dimensions of a length—is quite large in comparison with the thickness of the material.

32-5 The index of a mixture

There is another prediction of our theory of the index of refraction that we can check against experiment. Suppose we consider a mixture of two materials. The index of the mixture is not the average of the two indexes, but should be given in terms of the sum of the two polarizabilities, as in Eq. (32.34). If we ask about the index of, say, a sugar solution, the total polarizability is the sum of the polarizability of the water and that of the sugar. Each must, of course, be calculated using for N the number per unit volume of the molecules of the particular kind. In other words, if a given solution has N_1 molecules of water, whose polarizability is α_1 , and N_2 molecules of sucrose ($\text{C}_{12}\text{H}_{22}\text{O}_{11}$), whose polarizability is α_2 , we should have that

$$3 \left(\frac{n^2 - 1}{n^2 + 2} \right) = N_1 \alpha_1 + N_2 \alpha_2. \quad (32.37)$$

We can use this formula to test our theory against experiment by measuring the index for various concentrations of sucrose in water. We are making several assumptions here, however. Our formula assumes that there is no chemical action when the sucrose is dissolved and that the disturbances to the individual atomic

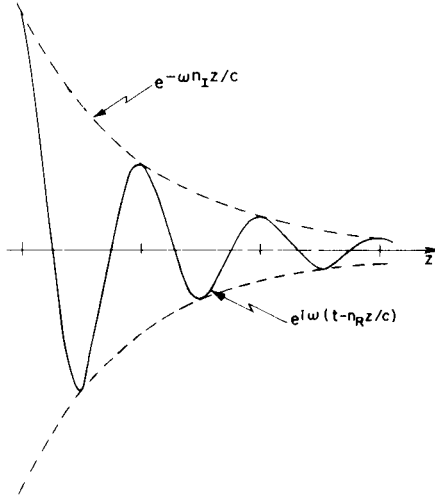


Fig. 32-1. A graph of E_x for some instant t , if $n_I \approx n_R/2\pi$.

Table 32-2

Refractive index of sucrose solutions, and comparison with predictions of Eq. (32.37).

Data from Handbook

A	B	C	D	E	F	G	H	J
Fraction of sucrose by weight	density (gm/cm ³)	n at 20°C	Moles of sucrose ^d per liter, N_2/N_0	Moles of water ^e per liter, N_1/N_0	$3 \left(\frac{n^2 - 1}{n^2 + 2} \right)$	$N_1\alpha_1$	$N_2\alpha_2$	$N_0\alpha_2$ (gm./liter)
0 ^a	0.9982	1.333	0	55.5	0.617	0.617	0	—
0.30	1.1270	1.3811	0.970	43.8	0.698	0.487	0.211	0.213
0.50	1.2296	1.4200	1.798	34.15	0.759	0.379	0.380	0.211
0.85	1.4454	1.5033	3.59	12.02	0.886	0.1335	0.752	0.210
1.00 ^b	1.588	1.5577 ^c	4.64	0	0.960	0	0.960	0.207

^a pure water^c average (see text)^e molecular weight of water = 18^b sugar crystals^d molecular weight of sucrose = 342

oscillators are not too different for various concentrations. So our result is certainly only approximate. Anyway, let's see how good it is.

We have picked the example of a sugar solution because there is a good table of measurements of the index of refraction in the *Handbook of Chemistry and Physics* and also because sugar is a molecular crystal that goes into solution without ionizing or otherwise changing its chemical state.

We give in the first three columns of Table 32-2 the data from the handbook. Column A is the percent of sucrose by weight, column B is the measured density (gm/cm³), and column C is the measured index of refraction for light whose wavelength is 589.3 millimicrons. For pure sugar we have taken the measured index of sugar crystals. The crystals are not isotropic, so the measured index is different along different directions. The handbook gives three values:

$$n_1 = 1.5376, \quad n_2 = 1.5651, \quad n_3 = 1.5705.$$

We have taken the average.

Now we could try to compute n for each concentration, but we don't know what value to take for α_1 or α_2 . Let's test the theory this way: We will assume that the polarizability of water (α_1) is the same at all concentrations and compute the polarizability of sucrose by using the experiment of values for n and solving Eq. (38.27) for α_2 . If the theory is correct, we should get the same α_2 for all concentrations.

First, we need to know N_1 and N_2 : let's express them in terms of Avogadro's number, N_0 . Let's take one liter (1000 cm³) for our unit of volume. Then N_i/N_0 is the weight per liter divided by the gram-molecular weight. And the weight per liter is the density (multiplied by 1000 to get grams per liter) times the fractional weight of either the sucrose or the water. In this way, we get N_2/N_0 and N_1/N_0 as in columns D and E of the table.

In column F we have computed $3(n^2 - 1)/(n^2 + 2)$ from the experimental values of n in column C. For pure water, $3(n^2 - 1)/(n^2 + 2)$ is 0.617, which is equal to just $N_1\alpha_1$. We can then fill in the rest of Column G, since for each row row G/E may be in the same ratio—namely, 0.617:55.5. Subtracting column G from column F, we get the contribution $N_2\alpha_2$ of the sucrose, shown in column H. Dividing these entries by the values of N_2/N_0 in column D, we get the value of $N_0\alpha_2$ shown in column J.

From our theory we would expect all the values of $N_0\alpha_2$ to be the same. They are not exactly equal, but pretty close. We can conclude that our ideas are fairly correct. Even more, we find that the polarizability of the sugar molecule doesn't seem to depend much on its surroundings—its polarizability is nearly the same in a dilute solution as it is in the crystal.

32-6 Waves in metals

The theory we have worked out in this chapter for solid materials can also be applied to good conductors, like metals, with very little modification. In metals some of the electrons have no binding force holding them to any particular atom; it is these “free” electrons which are responsible for the conductivity. There are other electrons which are bound, and the theory above is directly applicable to them. Their influence, however, is usually swamped by the effects of the conduction electrons. We will consider now only the effects of the free electrons.

If there is no restoring force on an electron—but still some resistance to its motion—its equation of motion differs from Eq. (32.1) only because the term in $\omega_0^2 x$ is lacking. So all we have to do is set $\omega_0^2 = 0$ in the rest of our derivations—except that there is one more difference. The reason that we had to distinguish between the average field and the local field in a dielectric is that in an insulator each of the dipoles is fixed in position, so that it has a definite relationship to the position of the others. But because the conduction electrons in a metal move around all over the place, the field on them *on the average* is just the average field E . So the correction we made to Eq. (32.5) by using Eq. (32.28) should *not* be made for conduction electrons. Therefore the formula for the index of refraction for metals should look like Eq. (32.27), except with ω_0 set equal to zero, namely,

$$n^2 = 1 + \frac{Nq_e^2}{m\epsilon_0} \frac{1}{-\omega^2 + i\gamma\omega}. \quad (32.38)$$

This is only the contribution from the conduction electrons, which we will assume is the major term for metals.

Now we even know how to find what value to use for γ , because it is related to the conductivity of the metal. In Chapter 43 of Volume I we discussed how the conductivity of a metal comes from the diffusion of the free electrons through the crystal. The electrons go on a jagged path from one scattering to the next, and between scatterings they move freely except for an acceleration due to any average electric field (as shown in Fig. 32-2). We found in Chapter 43 of Volume I that the average drift velocity is just the acceleration times the average time τ between collisions. The acceleration is $q_e E/m$, so

$$v_{\text{drift}} = \frac{q_e E}{m} \tau. \quad (32.39)$$

This formula assumed that E was constant, so that v_{drift} was a steady velocity. Since there is no average acceleration, the drag force is equal to the applied force. We have defined γ by saying that $\gamma m v$ is the drag force [see Eq. (32.1)], which is $q_e E$; therefore we have that

$$\gamma = \frac{1}{\tau}. \quad (32.40)$$

Although we cannot easily measure τ directly, we can determine it by measuring the conductivity of the metal. It is found experimentally that an electric field E in a metal produces a current with the density j proportional to E (for isotropic materials):

$$j = \sigma E.$$

The proportionality constant σ is called the *conductivity*. This is just what we expect from Eq. (32.39) if we set

$$j = Nq_e v_{\text{drift}}.$$

Then

$$\sigma = \frac{Nq_e^2}{m} \tau. \quad (32.41)$$

So τ —and therefore γ —can be related to the observed electrical conductivity. Using Eqs. (32.40) and (32.41), we can rewrite our formula for the index, Eq.

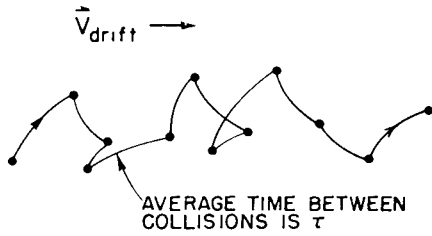


Fig. 32-2. The motion of a free electron.

(32.38), in the following form:

$$n^2 = 1 + \frac{\sigma/\epsilon_0}{i\omega(1 + i\omega\tau)}, \quad (32.42)$$

where

$$\tau = \frac{1}{\gamma} = \frac{m\sigma}{Nq_e^2}. \quad (32.43)$$

This is a convenient formula for the index of refraction of metals.

32-7 Low-frequency and high-frequency approximations; the skin depth and the plasma frequency

Our result, Eq. (32.42), for the index of refraction for metals predicts quite different characteristics for wave propagation at different frequencies. Let's first see what happens at very *low* frequencies. If ω is small enough, we can approximate Eq. (32.42) by

$$n^2 = -i \frac{\sigma}{\epsilon_0\omega}. \quad (32.44)$$

Now, as you can check by taking the square,*

$$\sqrt{-i} = \frac{1-i}{\sqrt{2}};$$

so for low frequencies,

$$n = \sqrt{\sigma/2\epsilon_0\omega} (1 - i). \quad (32.45)$$

The real and imaginary parts of n have the same magnitude. With such a large imaginary part to n , the wave is rapidly attenuated in the metal. Referring to Eq. (32.36), the amplitude of a wave going in the z -direction decreases as

$$\exp[-\sqrt{\sigma\omega/2\epsilon_0c^2} z]. \quad (32.46)$$

Let's write this as

$$e^{-z/\delta}, \quad (32.47)$$

where δ is then the distance in which the wave amplitude decreases by the factor $e^{-1} = 1/2.72$ —or roughly one-third. The amplitude of such a wave as a function of z is shown in Fig. 32-3. Since electromagnetic waves will penetrate into a metal only this distance, δ is called the *skin depth*. It is given by

$$\delta = \sqrt{2\epsilon_0c^2/\sigma\omega}. \quad (32.48)$$

Now what do we mean by “low” frequencies? Looking at Eq. (32.42), we see that it can be approximated by Eq. (32.44) only if $\omega\tau$ is much less than one and if $\omega\epsilon_0/\sigma$ is also much less than one—that is, our low-frequency approximation applies when

$$\omega \ll \frac{1}{\tau}$$

and

$$\omega \ll \frac{\sigma}{\epsilon_0}. \quad (32.49)$$

Let's see what frequencies these correspond to for a typical metal like copper. We compute τ by using Eq. (32.43), and σ/ϵ_0 , by using the measured conductivity. We take the following data from a handbook:

$$\begin{aligned} \sigma &= 5.76 \times 10^7 \text{ (ohm-meter)}^{-1}, \\ \text{atomic weight} &= 63.5 \text{ grams}, \\ \text{density} &= 8.9 \text{ grams/cm}^3, \\ \text{Avogadro's number} &= 6.02 \times 10^{23} \text{ (gram atomic weight)}^{-1}. \end{aligned}$$

* Or writing $-i = e^{-i\pi/2}$; $\sqrt{-i} = e^{-i\pi/4} = \cos \pi/4 - i \sin \pi/4$, which gives the same result.

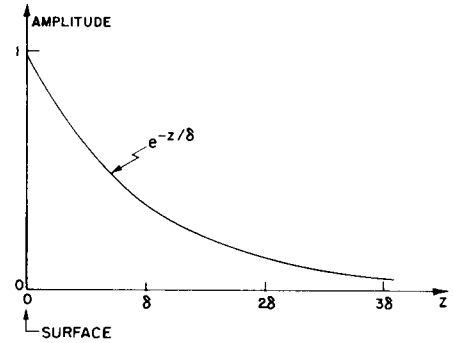


Fig. 32-3. The amplitude of a transverse electromagnetic wave as a function of distance into a metal.

If we assume that there is one free electron per atom, then the number of electrons per cubic meter is

$$N = 8.5 \times 10^{28} \text{ meter}^{-3}.$$

Using

$$q_e = 1.6 \times 10^{-19} \text{ coulomb},$$

$$\epsilon_0 = 8.85 \times 10^{-12} \text{ farad-meter}^{-1},$$

$$m = 9.11 \times 10^{-31} \text{ kgm},$$

we get

$$\tau = 2.4 \times 10^{-14} \text{ sec},$$

$$\frac{1}{\tau} = 4.1 \times 10^{13} \text{ sec}^{-1},$$

$$\frac{\sigma}{\epsilon_0} = 6.5 \times 10^{18} \text{ sec}^{-1}.$$

So for frequencies less than about 10^{12} cycles per second, copper will have the “low-frequency” behavior we describe (that means for waves whose free-space wavelength is longer than 0.3 millimeters—*very* short radio waves!).

For these waves, the skin depth in copper is

$$\delta = \sqrt{\frac{0.028 \text{ m}^2\text{-sec}^{-1}}{\omega}}.$$

For microwaves of 10,000 megacycles per second (3-cm waves)

$$\delta = 6.7 \times 10^{-4} \text{ cm}.$$

The wave penetrates a very small distance.

We can see from this why in studying cavities (or waveguides) we needed to worry only about the fields inside the cavity, and not in the metal or outside the cavity. Also, we see why the losses in a cavity are reduced by a thin plating of silver or gold. The losses come from the current, which are appreciable only in a thin layer equal to the skin depth.

Suppose we look now at the index of a metal like copper at high frequencies. For very high frequencies $\omega\tau$ is much greater than one, and Eq. (32.42) is well approximated by

$$n^2 = 1 - \frac{\sigma}{\epsilon_0\omega^2\tau}. \quad (32.50)$$

For waves of high frequencies the index of a metal becomes real—and less than one! This is also evident from Eq. (32.38) if the dissipation term with γ is neglected, as can be done for very large ω . Equation (32.38) gives

$$n^2 = 1 - \frac{Nq_e^2}{m\epsilon_0\omega^2}, \quad (32.51)$$

which is, of course, the same as Eq. (32.50). We have seen before the quantity $Nq_e^2/m\epsilon_0$, which we called the square of the plasma frequency (Section 7-3):

$$\omega_p^2 = \frac{Nq_e^2}{\epsilon_0 m},$$

so we can write Eq. (32.50) or Eq. (32.51) as

$$n^2 = 1 - \left(\frac{\omega_p}{\omega}\right)^2.$$

The plasma frequency is a kind of “critical” frequency.

For $\omega < \omega_p$ the index of a metal has an imaginary part, and waves are attenuated; but for $\omega \gg \omega_p$ the index is real, and the metal becomes transparent. You know, of course, that metals are reasonably transparent to x-rays. But some metals are even transparent in the ultraviolet. In Table 32-3 we give for

several metals the experimental observed wavelength at which they begin to become transparent. In the second column we give the calculated critical wavelength $\lambda_p = 2\pi c/\omega_p$. Considering that the experimental wavelength is not too well defined, the fit of the theory is fairly good.

You may wonder why the plasma frequency ω_p should have anything to do with the propagation of electromagnetic waves in metals. The plasma frequency came up in Chapter 7 as the natural frequency of *density* oscillations of the free electrons. (A clump of electrons is repelled by electric forces, and the inertia of the electrons leads to an oscillation of density.) So *longitudinal* plasma waves are resonant at ω_p . But we are now talking about *transverse* electromagnetic waves, and we have found that transverse waves are absorbed for frequencies below ω_p . (It's an interesting and *not* accidental coincidence.)

Although we have been talking about wave propagation in metals, you appreciate by this time the universality of the phenomena of physics—that it doesn't make any difference whether the free electrons are in a metal or whether they are in the plasma of the ionosphere of the earth, or in the atmosphere of a star. To understand radio propagation in the ionosphere, we can use the same expressions—using, of course, the proper values for N and τ . We can see now why long radio waves are absorbed or reflected by the ionosphere, whereas short waves go right through. (Short waves must be used for communication with satellites.)

We have talked about the high- and low-frequency extremes for wave propagation in metals. For the in-between frequencies the full-blown formula of Eq. (32.42) must be used. In general, the index will have real and imaginary parts; the wave is attenuated as it propagates into the metal. For very thin layers, metals are somewhat transparent even at optical frequencies. As an example, special goggles for people who work around high-temperature furnaces are made by evaporating a thin layer of gold on glass. The visible light is transmitted fairly well—with a strong green tinge—but the infrared is strongly absorbed.

Finally, it cannot have escaped the reader that many of these formulas resemble in some ways those for the dielectric constant κ discussed in Chapter 10. The dielectric constant κ measures the response of the material to a constant field, that is, for $\omega = 0$. If you look carefully at the definition of n and κ you see that κ is simply the limit of n^2 as $\omega \rightarrow 0$. Indeed, placing $\omega = 0$ and $n^2 = \kappa$ in equations of this chapter will reproduce the equations of the theory of the dielectric constant of Chapter 11.

Table 32-3*

Wavelengths below which the metal becomes transparent

Metal	λ (experimental)	$\lambda_p = 2\pi c/\omega_p$
Li	1550 Å	1550 Å
Na	2100	2090
K	3150	2870
Rb	3400	3220

* From: C. Kittel, *Introduction to Solid State Physics*, John Wiley and Sons, Inc., New York, 2nd ed., 1956, p. 266.

Reflection from Surfaces

33-1 Reflection and refraction of light

The subject of this chapter is the reflection and refraction of light—or electromagnetic waves in general—at surfaces. We have already discussed the laws of reflection and refraction in Chapter 35 of Volume I. Here's what we found out there:

1. The angle of reflection is equal to the angle of incidence. With the angles defined as shown in Fig. 33-1,

$$\theta_r = \theta_i. \quad (33.1)$$

2. The product $n \sin \theta$ is the same for the incident and transmitted beams (Snell's law).

$$n_1 \sin \theta_i = n_2 \sin \theta_t. \quad (33.2)$$

3. The intensity of the reflected light depends on the angle of incidence and also on the direction of polarization. For E perpendicular to the plane of incidence, the reflection coefficient R_{\perp} is

$$R_{\perp} = \frac{I_r}{I_i} = \frac{\sin^2 (\theta_i - \theta_t)}{\sin^2 (\theta_i + \theta_t)} \quad (33.3)$$

For E parallel to the plane of incidence, the reflection coefficient R_{\parallel} is

$$R_{\parallel} = \frac{I_r}{I_i} = \frac{\tan^2 (\theta_i - \theta_t)}{\tan^2 (\theta_i + \theta_t)}. \quad (33.4)$$

4. For normal incidence (any polarization, of course!),

$$\frac{I_r}{I_i} = \left(\frac{n_2 - n_1}{n_2 + n_1} \right)^2. \quad (33.5)$$

(Earlier, we used i for the incident angle and r for the refracted angle. Since we can't use r for both "refracted" and "reflected" angles, we are now using θ_i = incident angle, θ_r = reflected angle, and θ_t = transmitted angle.)

Our earlier discussion is really about as far as anyone would normally need to go with the subject, but we are going to do it all over again a different way. Why? One reason is that we assumed before that the indexes were real (no absorption in the materials). But another reason is that you should know how to deal with what happens to waves at surfaces from the point of view of Maxwell's equations. We'll get the same answers as before, but now from a straightforward solution of the wave problem, rather than by some clever arguments.

We want to emphasize that the amplitude of a surface reflection is not a property of the *material*, as is the index of refraction. It is a "surface property," one that depends precisely on how the surface is made. A thin layer of extraneous junk on the surface between two materials of indices n_1 and n_2 will usually change the reflection. (There are all kinds of possibilities of interference here—like the colors of oil films. Suitable thickness can even reduce the reflected amplitude to zero for a given frequency; that's how coated lenses are made.) The formulas we will derive are correct only if the change of index is sudden—within a distance very small compared with one wavelength. For light, the wavelength is about 5000 Å, so by a "smooth" surface we mean one in which the conditions change in

33-1 Reflection and refraction of light

33-2 Waves in dense materials

33-3 The boundary conditions

33-4 The reflected and transmitted waves

33-5 Reflection from metals

33-6 Total internal reflection

Review. Chapter 35, Vol. I, Polarization

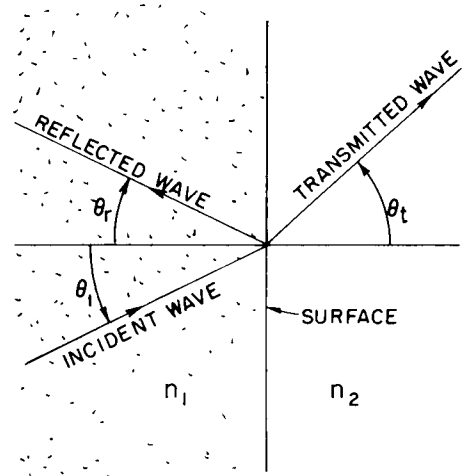


Fig. 33-1. Reflection and refraction of light waves at a surface. (The wave directions are normal to the wave crests.)

going a distance of only a few atoms (or a few angstroms). Our equations will work for light for highly polished surfaces. In general, if the index changes gradually over a distance of several wavelengths, there is very little reflection at all.

33-2 Waves in dense materials

First, we remind you about the convenient way of describing a sinusoidal plane wave we used in Chapter 36 of Volume I. Any field *component* in the wave (we use E as an example) can be written in the form

$$E = E_0 e^{i(\omega t - \mathbf{k} \cdot \mathbf{r})}, \quad (33.6)$$

where E represents the amplitude at the point \mathbf{r} (from the origin) at the time t . The vector \mathbf{k} points in the direction the wave is travelling, and its magnitude $|\mathbf{k}| = k = 2\pi/\lambda$ is the wave number. The phase velocity of the wave is $v_{ph} = \omega/k$, for a light wave in a material of index n , $v_{ph} = c/n$, so

$$k = \frac{\omega n}{c}. \quad (33.7)$$

Suppose \mathbf{k} is in the z -direction, then $\mathbf{k} \cdot \mathbf{r}$ is just kz , as we have often used it. For \mathbf{k} in any other direction, we should replace z by r_k , the distance from the origin in the \mathbf{k} -direction; that is, we should replace kz by kr_k , which is just $\mathbf{k} \cdot \mathbf{r}$. (See Fig. 33-2.) So Eq. (33.6) is a convenient representation of a wave in any direction.

We must remember, of course, that

$$\mathbf{k} \cdot \mathbf{r} = k_x x + k_y y + k_z z,$$

where k_x , k_y , and k_z are the components of \mathbf{k} along the three axes. In fact, we pointed out once that (ω, k_x, k_y, k_z) is a four-vector, and that its scalar product with (t, x, y, z) is an invariant. So the *phase* of a wave is an invariant, and Eq. (33.6) could be written

$$E = E_0 e^{ik_\mu x_\mu}.$$

But we don't need to be that fancy now.

For a sinusoidal E , as in Eq. (33.6), $\partial E/\partial t$ is the same as $i\omega E$, and $\partial E/\partial x$ is $-ik_x E$, and so on for the other components. You can see why it is very convenient to use the form in Eq. (33.6) when working with differential equations—differentiations are replaced by multiplications. One further useful point: The operation $\nabla = (\partial/\partial x, \partial/\partial y, \partial/\partial z)$ gets replaced by the three multiplications $(-ik_x, -ik_y, -ik_z)$. But these three factors transform as the components of the vector \mathbf{k} , so the operator ∇ gets replaced by multiplication with $-i\mathbf{k}$:

$$\begin{aligned} \frac{\partial}{\partial t} &\rightarrow i\omega, \\ \nabla &\rightarrow -i\mathbf{k}. \end{aligned} \quad (33.8)$$

This remains true for any ∇ operation—whether it is the gradient, or the divergence, or the curl. For instance, the z -component of $\nabla \times \mathbf{E}$ is

$$\frac{\partial E_y}{\partial x} - \frac{\partial E_x}{\partial y}.$$

If both E_y and E_x vary as $e^{-i\mathbf{k} \cdot \mathbf{r}}$, then we get

$$-ik_x E_y + ik_y E_x,$$

which is, you see, the z -component of $-i\mathbf{k} \times \mathbf{E}$.

So we have the very useful general fact that whenever you have to take the gradient of a vector that varies as a wave in three dimensions (they are an important part of physics), you can always take the derivations quickly and almost without thinking by remembering that the operation ∇ is equivalent to multiplication by $-i\mathbf{k}$.

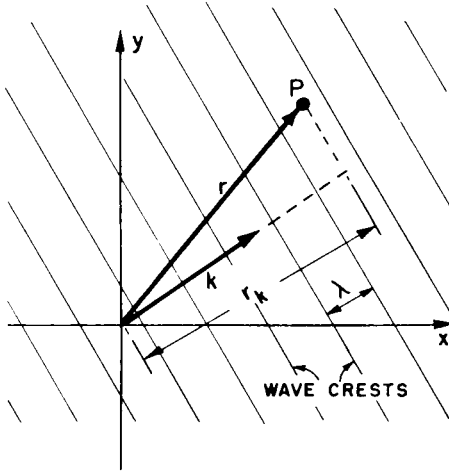


Fig. 33-2. For a wave moving in the direction \mathbf{k} , the phase at any point P is $(\omega t - \mathbf{k} \cdot \mathbf{r})$.

For instance, the Faraday equation

$$\nabla \times \mathbf{E} = -\frac{\partial \mathbf{B}}{\partial t}$$

becomes for a wave

$$-ik \times \mathbf{E} = -i\omega \mathbf{B}.$$

This tells us that

$$\mathbf{B} = \frac{\mathbf{k} \times \mathbf{E}}{\omega}, \quad (33.9)$$

which corresponds to the result we found earlier for waves in free space—that \mathbf{B} , in a wave, is at right angles to \mathbf{E} and to the wave direction. (In free space, $\omega/k = c$.) You can remember the sign in Eq. (33.9) from the fact that \mathbf{k} is in the direction of Poynting's vector $\mathbf{S} = \epsilon_0 c^2 \mathbf{E} \times \mathbf{B}$.

If you use the same rule with the other Maxwell equations, you get again the results of the last chapter and, in particular, that

$$\mathbf{k} \cdot \mathbf{k} = k^2 = \frac{\omega^2 n^2}{c^2}. \quad (33.10)$$

But since we know that, we won't do it again.

If you want to entertain yourself, you can try the following terrifying problem that was the ultimate test for graduate students back in 1890: solve Maxwell's equations for plane waves in an *anisotropic* crystal, that is, when the polarization \mathbf{P} is related to the electric field \mathbf{E} by a tensor of polarizability. You should, of course, choose your axes along the principal axes of the tensor, so that the relations are simplest (then $P_x = \alpha_a E_x$, $P_y = \alpha_b E_y$, and $P_z = \alpha_c E_z$), but let the waves have an arbitrary direction and polarization. You should be able to find the relations between \mathbf{E} and \mathbf{B} , and how \mathbf{k} varies with direction and wave polarization. Then you will understand the optics of an anisotropic crystal. It would be best to start with the simpler case of a birefringent crystal—like calcite—for which two of the polarizabilities are equal (say, $\alpha_b = \alpha_c$), and see if you can understand why you see double when you look through such a crystal. If you can do that, then try the hardest case, in which all three α 's are different. Then you will know whether you are up to the level of a graduate student of 1890. In this chapter, however, we will consider only isotropic substances

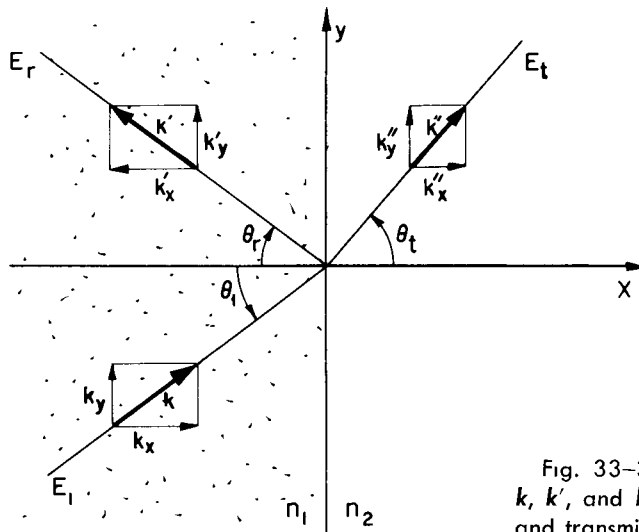


Fig. 33-3. The propagation vectors \mathbf{k} , \mathbf{k}' , and \mathbf{k}'' for the incident, reflected, and transmitted waves.

We know from experience that when a plane wave arrives at the boundary between two different materials—say, air and glass, or water and oil—there is a wave reflected and a wave transmitted. Suppose we assume no more than that and see what we can work out. We choose our axes with the yz -plane in the surface and the xy -plane perpendicular to the incident wave surfaces, as shown in Fig. 33-3.

The electric vector of the incident wave can then be written as

$$\mathbf{E}_i = \mathbf{E}_0 e^{i(\omega t - \mathbf{k} \cdot \mathbf{r})}. \quad (33.11)$$

Since \mathbf{k} is perpendicular to the z -axis,

$$\mathbf{k} \cdot \mathbf{r} = k_x x + k_y y. \quad (33.12)$$

We write the reflected wave as

$$\mathbf{E}_r = \mathbf{E}'_0 e^{i(\omega' t - \mathbf{k}' \cdot \mathbf{r})}, \quad (33.13)$$

so that its frequency is ω' , its wave number is \mathbf{k}' , and its amplitude is \mathbf{E}'_0 . (We know, of course, that the frequency is the same and the magnitude of \mathbf{k} is the same as for the incident wave, but we are not going to assume even that. We will let it come out of the mathematical machinery.) Finally, we write for the transmitted wave,

$$\mathbf{E}_t = \mathbf{E}''_0 e^{i(\omega'' t - \mathbf{k}'' \cdot \mathbf{r})}. \quad (33.14)$$

We know that one of Maxwell's equations gives Eq (33.9), so for each of the waves we have

$$\mathbf{B}_i = \frac{\mathbf{k} \times \mathbf{E}_i}{\omega}, \quad \mathbf{B}_r = \frac{\mathbf{k}' \times \mathbf{E}_r}{\omega'}, \quad \mathbf{B}_t = \frac{\mathbf{k}'' \times \mathbf{E}_t}{\omega''}. \quad (33.15)$$

Also, if we call the indexes of the two media n_1 and n_2 , we have from Eq. (33.10)

$$k^2 = k_x^2 + k_y^2 = \frac{\omega^2 n_1^2}{c^2}. \quad (33.16)$$

Since the reflected wave is in the same material, then

$$k'^2 = \frac{\omega'^2 n_1^2}{c^2}, \quad (33.17)$$

whereas for the transmitted wave,

$$k''^2 = \frac{\omega''^2 n_2^2}{c^2}. \quad (33.18)$$

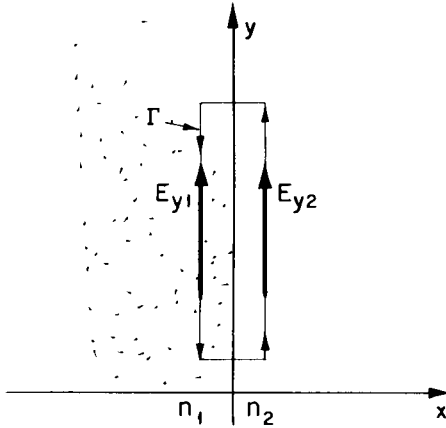


Fig. 33-4. A boundary condition $E_{y2} = E_{y1}$ is obtained from $\oint_{\Gamma} \mathbf{E} \cdot d\mathbf{s} = 0$.

33-3 The boundary conditions

All we have done so far is to describe the three waves; our problem now is to work out the parameters of the reflected and transmitted waves in terms of those of the incident wave. How can we do that? The three waves we have described satisfy Maxwell's equations in the uniform material, but Maxwell's equations must also be satisfied *at* the boundary between the two different materials. So we must now look at what happens right *at* the boundary. We will find that Maxwell's equations demand that the three waves fit together in a certain way.

As an example of what we mean, the y -component of the electric field \mathbf{E} must be the *same* on both sides of the boundary. This is required by Faraday's law,

$$\nabla \times \mathbf{E} = - \frac{\partial \mathbf{B}}{\partial t}, \quad (33.19)$$

as we can see in the following way. Consider a little rectangular loop Γ which straddles the boundary, as shown in Fig 33-4. Equation (33.19) says that the line integral of \mathbf{E} around Γ is equal to the rate of change of the flux of \mathbf{B} through the loop:

$$\oint_{\Gamma} \mathbf{E} \cdot d\mathbf{s} = - \frac{\partial}{\partial t} \int \mathbf{B} \cdot \mathbf{n} \, da.$$

Now imagine that the rectangle is very narrow, so that the loop encloses an infinitesimal area. If \mathbf{B} remains finite (and there's no reason it should be infinite at the boundary!) the flux through the area is zero. So the line integral of \mathbf{E} must

be zero. If E_{y1} and E_{y2} are the components of the field on the two sides of the boundary and if the length of the rectangle is l , we have

$$E_{y1}l - E_{y2}l = 0$$

or

$$E_{y1} = E_{y2}, \quad (33.20)$$

as we have said. This gives us one relation among the fields of the three waves.

The procedure of working out the consequences of Maxwell's equations at the boundary is called "determining the boundary conditions." Ordinarily, it is done by finding as many equations like Eq. (33.20) as one can, by making arguments about little rectangles like Γ in Fig. 33-4, or by using little gaussian surfaces that straddle the boundary. Although that is a perfectly good way of proceeding, it gives the impression that the problem of dealing with a boundary is different for every different physical problem.

For example, in a problem of heat flow across a boundary, how are the temperatures on the two sides related? Well, you could argue, for one thing, that the heat flow *to* the boundary from one side would have to equal the flow *away* from the other side. It is usually possible, and generally quite useful, to work out the boundary conditions by making such physical arguments. There may be times, however, when in working on some problem you have only some equations, and you may not see right away what physical arguments to use. So although we are at the moment interested only in an electromagnetic problem, where we *can* make the physical arguments, we want to show you a method that can be used for any problem—a *general* way of finding what happens at a boundary directly from the differential equations.

We begin by writing all the Maxwell equations for a dielectric—and this time we are very specific and write out explicitly all the components:

$$\begin{aligned} \nabla \cdot \mathbf{E} &= -\frac{\nabla \cdot \mathbf{P}}{\epsilon_0} \\ \epsilon_0 \left(\frac{\partial E_x}{\partial x} + \frac{\partial E_y}{\partial y} + \frac{\partial E_z}{\partial z} \right) &= - \left(\frac{\partial P_x}{\partial x} + \frac{\partial P_y}{\partial y} + \frac{\partial P_z}{\partial z} \right) \end{aligned} \quad (33.21)$$

$$\begin{aligned} \nabla \times \mathbf{E} &= -\frac{\partial \mathbf{B}}{\partial t} \\ \frac{\partial E_z}{\partial y} - \frac{\partial E_y}{\partial z} &= -\frac{\partial B_x}{\partial t} \end{aligned} \quad (33.22a)$$

$$\frac{\partial E_x}{\partial z} - \frac{\partial E_z}{\partial x} = -\frac{\partial B_y}{\partial t} \quad (33.22b)$$

$$\frac{\partial E_y}{\partial x} - \frac{\partial E_x}{\partial y} = -\frac{\partial B_z}{\partial t} \quad (33.22c)$$

$$\begin{aligned} \nabla \cdot \mathbf{B} &= 0 \\ \frac{\partial B_x}{\partial x} + \frac{\partial B_y}{\partial y} + \frac{\partial B_z}{\partial z} &= 0 \end{aligned} \quad (33.23)$$

$$\begin{aligned} c^2 \nabla \times \mathbf{B} &= \frac{1}{\epsilon_0} \frac{\partial \mathbf{P}}{\partial t} + \frac{\partial \mathbf{E}}{\partial t} \\ c^2 \left(\frac{\partial B_z}{\partial y} - \frac{\partial B_y}{\partial z} \right) &= \frac{1}{\epsilon_0} \frac{\partial P_x}{\partial t} + \frac{\partial E_x}{\partial t} \end{aligned} \quad (33.24a)$$

$$c^2 \left(\frac{\partial B_x}{\partial z} - \frac{\partial B_z}{\partial x} \right) = \frac{1}{\epsilon_0} \frac{\partial P_y}{\partial t} + \frac{\partial E_y}{\partial t} \quad (33.24b)$$

$$c^2 \left(\frac{\partial B_y}{\partial x} - \frac{\partial B_x}{\partial y} \right) = \frac{1}{\epsilon_0} \frac{\partial P_z}{\partial t} + \frac{\partial E_z}{\partial t} \quad (33.24c)$$

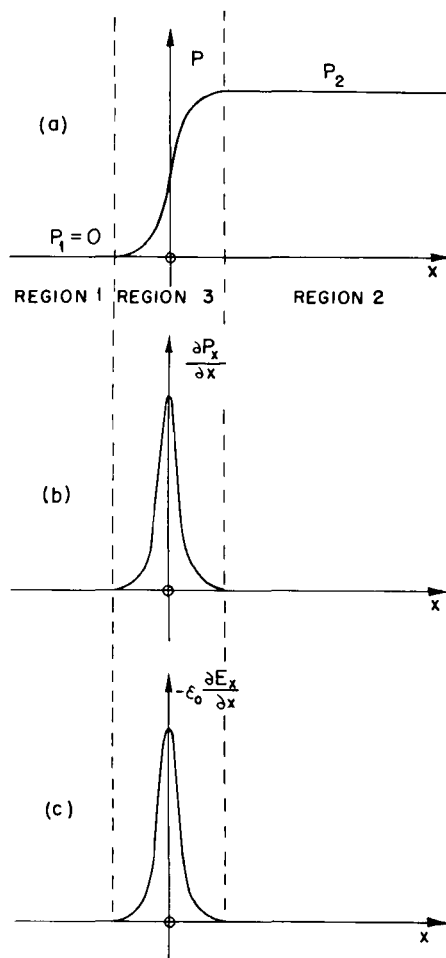


Fig. 33-5. The fields in the transition region (3) between two different materials in regions (1) and (2).

Now these equations must all hold in region 1 (to the left of the boundary) and in region 2 (to the right of the boundary). We have already written the solutions in regions 1 and 2. Finally, they must also be satisfied *in* the boundary, which we can call region 3. Although we usually think of the boundary as being sharply discontinuous, in reality it is not. The physical properties change very rapidly but not infinitely fast. In any case, we can imagine that there is a very rapid, but *continuous*, transition of the index between region 1 and 2, in a short distance we can call region 3. Also, any field quantity like P_x , or E_y , etc., will make a similar kind of transition in region 3. In this region, the differential equations must still be satisfied, and it is by following the differential equations in this region that we can arrive at the needed “boundary conditions.”

For instance, suppose that we have a boundary between vacuum (region 1) and glass (region 2). There is nothing to polarize in the vacuum, so $P_1 = 0$. Let's say there is some polarization P_2 in the glass. Between the vacuum and the glass there is a smooth, but rapid, transition. If we look at any component of P , say P_x , it might vary as drawn in Fig. 33-5(a). Suppose now we take the first of our equations, Eq. (33.21). It involves derivatives of the components of P with respect to x , y , and z . The y - and z -derivatives are not interesting; nothing spectacular is happening in those directions. But the x -derivative of P_x will have some very large values in region 3, because of the tremendous slope of P_x . The derivative $\partial P_x / \partial x$ will have a sharp spike at the boundary, as shown in Fig. 33-5(b). If we imagine squashing the boundary to an even thinner layer, the spike would get much higher. If the boundary is really sharp for the waves we are interested in, the magnitude of $\partial P_x / \partial x$ in region 3 will be much, much greater than any contributions we might have from the variation of P in the wave away from the boundary—so we ignore any variations other than those due to the boundary.

Now how can Eq. (33.21) be satisfied if there is a whopping big spike on the right-hand side? Only if there is an equally whopping big spike on the other side. Something on the left-hand side must also be big. The only candidate is $\partial E_x / \partial x$, because the variations with y and z are only those small effects in the wave we just mentioned. So $-\epsilon_0(\partial E_x / \partial x)$ must be as drawn in Fig. 33-5(c)—just a copy of $\partial P_x / \partial x$. We have that

$$\epsilon_0 \frac{\partial E_x}{\partial x} = - \frac{\partial P_x}{\partial x}.$$

If we integrate this equation with respect to x across region 3, we conclude that

$$\epsilon_0(E_{x2} - E_{x1}) = -(P_{x2} - P_{x1}). \quad (33.25)$$

In other words, the jump in $\epsilon_0 E_x$ in going from region 1 to region 2 must be equal to the jump in $-P_x$.

We can rewrite Eq. (33.25) as

$$\epsilon_0 E_{x2} + P_{x2} = \epsilon_0 E_{x1} + P_{x1}, \quad (33.26)$$

which says that the quantity $(\epsilon_0 E_x + P_x)$ has equal values in region 2 and region 1. People say: the quantity $(\epsilon_0 E_x + P_x)$ is *continuous* across the boundary. We have, in this way, one of our boundary conditions.

Although we took as an illustration the case in which P_1 was zero because region 1 was a vacuum, it is clear that the same argument applies for any two materials in the two regions, so Eq. (33.26) is true in general.

Let's now go through the rest of Maxwell's equations and see what each of them tells us. We take next Eq. (33.22a). There are no x -derivatives, so it doesn't tell us anything. (Remember that the fields *themselves* do not get especially large at the boundary; only the derivatives with respect to x can become so huge that they dominate the equation.) Next, we look at Eq. (33.22b). Ah! There is an x -derivative! We have $\partial E_z / \partial x$ on the left-hand side. Suppose it has a huge derivative. But wait a moment! There is nothing on the right-hand side to match it with; therefore E_z *cannot* have any jump in going from region 1 to region 2. [If it did, there would be a spike on the left of Eq. (33.22a) but none on the right,

and the equation would be false.] So we have a new condition:

$$E_{z2} = E_{z1}. \quad (33.27)$$

By the same argument, Eq. (33.22c) gives

$$E_{y2} = E_{y1}. \quad (33.28)$$

This last result is just what we got in Eq. (33.20) by a line integral argument.

We go on to Eq. (33.23). The only term that could have a spike is $\partial B_x / \partial x$. But there's nothing on the right to match it, so we conclude that

$$B_{x2} = B_{x1}. \quad (33.29)$$

On to the last of Maxwell's equations! Equation (33.24a) gives nothing, because there are no λ -derivatives. Equation (33.23b) has one, $-c^2 \partial B_z / \partial x$, but again, there is nothing to match it with. We get

$$B_{z2} = B_{z1}. \quad (33.30)$$

The last equation is quite similar, and gives

$$B_{y2} = B_{y1}. \quad (33.31)$$

The last three equations give us that $\mathbf{B}_2 = \mathbf{B}_1$. We want to emphasize, however, that we get this result only when the materials on both sides of the boundary are nonmagnetic—or rather, when we can neglect any magnetic effects of the materials. This can usually be done for most materials, except ferromagnetic ones. (We will treat the magnetic properties of materials in some later chapters.)

Our program has netted us the six relations between the fields in region 1 and those in region 2. We have put them all together in Table 33-1. We can now use them to match the waves in the two regions. We want to emphasize, however, that the idea we have just used will work in *any* physical situation in which you have differential equations and you want a solution that crosses a sharp boundary between two regions where some property changes. For our present purposes, we could have easily derived the same equations by using arguments about the fluxes and circulations at the boundary. (You might see whether you can get the same result that way.) But now you have seen a method that will work in case you ever get stuck and don't see any easy argument about the physics of what is happening at the boundary—you can just work with the equations.

Table 33-1

Boundary conditions at the surface of a dielectric

$$(\epsilon_0 \mathbf{E}_1 + \mathbf{P}_1)_x = (\epsilon_0 \mathbf{E}_2 + \mathbf{P}_2)_x$$

$$(\mathbf{E}_1)_y = (\mathbf{E}_2)_y$$

$$(\mathbf{E}_1)_z = (\mathbf{E}_2)_z$$

$$\mathbf{B}_1 = \mathbf{B}_2$$

(The surface is in the yz -plane)

33-4 The reflected and transmitted waves

Now we are ready to apply our boundary conditions to the waves we wrote down in Section 33-2. We had:

$$\mathbf{E}_i = E_0 e^{i(\omega t - k_x x - k_y y)}, \quad (33.32)$$

$$\mathbf{E}_r = E'_0 e^{i(\omega' t - k'_x x - k'_y y)}, \quad (33.33)$$

$$\mathbf{E}_t = E''_0 e^{i(\omega'' t - k''_x x - k''_y y)}, \quad (33.34)$$

$$\mathbf{B}_i = \frac{\mathbf{k} \times \mathbf{E}_i}{\omega}, \quad (33.35)$$

$$\mathbf{B}_r = \frac{\mathbf{k}' \times \mathbf{E}_r}{\omega'}, \quad (33.36)$$

$$\mathbf{B}_t = \frac{\mathbf{k}'' \times \mathbf{E}_t}{\omega''}. \quad (33.37)$$

We have one further bit of knowledge: \mathbf{E} is perpendicular to its propagation vector \mathbf{k} for each wave.

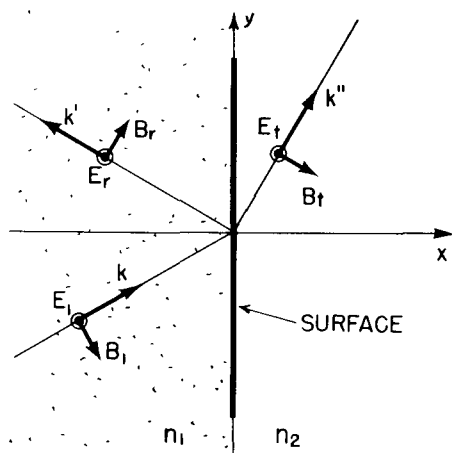


Fig. 33-6. Polarization of the reflected and transmitted waves when the E -field of the incident wave is perpendicular to the plane of incidence.

The results will depend on the direction of the E -vector (the “polarization”) of the incoming wave. The analysis is much simplified if we treat separately the case of an incident wave with its E -vector *parallel* to the “plane of incidence” (that is, the xy -plane) and the case of an incident wave with the E -vector *perpendicular* to the plane of incidence. A wave of any other polarization is just a linear combination of two such waves. In other words, the reflected and transmitted intensities are different for different polarizations, and it is easiest to pick the two simplest cases and treat them separately.

We will carry through the analysis for an incoming wave polarized perpendicular to the plane of incidence and then just give you the result for the other. We are cheating a little by taking the simplest case, but the principle is the same for both. So we take that E_i has only a z -component, and since all the E -vectors are in the same direction we can leave off the vector signs.

So long as both materials are isotropic, the induced oscillations of charges in the material will also be in the z -direction, and the E -field of the transmitted and radiated waves will have only z -components. So for all the waves, E_x and E_y and P_x and P_y are zero. The waves will have their E - and B -vectors as drawn in Fig. 33-6 (We are cutting a corner here on our original plan of getting everything from the equations. This result would also come out of the boundary conditions, but we can save a lot of algebra by using the physical argument. When you have some spare time, see if you can get the same result from the equations. It is clear that what we have said agrees with the equations; it is just that we have not shown that there are no *other* possibilities.)

Now our boundary conditions, Eqs. (33.26) through (33.31), give relations between the components of E and B in regions 1 and 2. For region 2 we have only the transmitted wave, but in region 1 we have *two* waves. Which one do we use? The fields in region 1 are, of course, the superposition of the fields of the incident and reflected waves. (Since each satisfies Maxwell’s equations, so does the sum.) So when we use the boundary conditions, we must use that

$$E_1 = E_i + E_r, \quad E_2 = E_t,$$

and similarly for the B ’s.

For the polarization we are considering, Eqs. (33.26) and (33.28) give us no new information; only Eq. (33.27) is useful. It says that

$$E_i + E_r = E_t$$

at the boundary, that is, for $x = 0$. So we have that

$$E_0 e^{i(\omega t - k_y y)} + E'_0 e^{i(\omega' t - k'_y y)} = E''_0 e^{i(\omega'' t - k''_y y)}, \quad (33.38)$$

which must be true for *all* t and for *all* y . Suppose we look first at $y = 0$. Then we have

$$E_0 e^{i\omega t} + E'_0 e^{i\omega' t} = E''_0 e^{i\omega'' t}$$

This equation says that two oscillating terms are equal to a third oscillation. That can happen only if all the oscillations have the same frequency. (It is impossible for three—or any number—of such terms with different frequencies to add to zero for all times.) So

$$\omega'' = \omega' = \omega. \quad (33.39)$$

As we knew all along, the frequencies of the reflected and transmitted waves are the same as that of the incident wave.

We should really have saved ourselves some trouble by putting that in at the beginning, but we wanted to show you that it can also be got out of the equations. When you are doing a real problem, it is usually the best thing to put everything you know into the works right at the start and save yourself a lot of trouble.

By definition, the *magnitude* of k is given by $k^2 = n^2 \omega^2 / c^2$, so we have also that

$$\frac{k''^2}{n_2^2} = \frac{k'^2}{n_1^2} = \frac{k^2}{n_1^2} \quad (33.40)$$

Now look at Eq. (33.38) for $t = 0$. Using again the same kind of argument we have just made, but this time based on the fact that the equation must hold for all values of y , we get that

$$k''_y = k'_y = k_y. \quad (33.41)$$

From Eq. (33.40), $k'^2 = k^2$, so

$$k'^2_x + k'^2_y = k^2_x + k^2_y.$$

Combining this with Eq. (33.41), we have that

$$k'^2_x = k^2_x,$$

or that $k'_x = \pm k_x$. The positive sign makes no sense; that would not give a *reflected* wave, but another *incident* wave, and we said at the start that we were solving the problem of only one incident wave. So we have

$$k'_x = -k_x. \quad (33.42)$$

The two equations (33.41) and (33.42) give us that the angle of reflection is equal to the angle of incidence, as we expected. (See Fig. 33-3) The reflected wave is

$$E_r = E'_0 e^{i(\omega t - k_x x + k_y y)}. \quad (33.43)$$

For the transmitted wave we already have that

$$k''_y = k_y,$$

and

$$\frac{k''^2}{n_2^2} = \frac{k^2}{n_1^2}; \quad (33.44)$$

so we can solve these to find k''_x . We get

$$k''^2_x = k''^2 - k''^2_y = \frac{n_2^2}{n_1^2} k^2 - k^2_y. \quad (33.45)$$

Suppose for a moment that n_1 and n_2 are real numbers (that the imaginary parts of the indexes are very small). Then all the k 's are also real numbers, and from Fig. 33-3 we find that

$$\frac{k_y}{k} = \sin \theta_i, \quad \frac{k''_y}{k''} = \sin \theta_t. \quad (33.46)$$

From (33.44) we get that

$$n_2 \sin \theta_t = n_1 \sin \theta_i, \quad (33.47)$$

which is Snell's law of refraction—again, something we already knew. If the indexes are not real, the wave numbers are complex, and we have to use Eq. (33.45). [We could still *define* the angles θ_i and θ_t by Eq. (33.46), and Snell's law, Eq. (33.47), would be true in general. But then the “angles” also are complex numbers, thereby losing their simple geometrical interpretation as angles. It is best then to describe the behavior of the waves by their complex k_x or k'_x values]

So far, we haven't found anything new. We have just had the simple-minded delight of getting some obvious answers from a complicated mathematical machinery. Now we are ready to find the amplitudes of the waves which we have not yet known. Using our results for the ω 's and k 's, the exponential factors in Eq. (33.38) can be cancelled, and we get

$$E_0 + E'_0 = E''_0. \quad (33.48)$$

Since both E'_0 and E''_0 are unknown, we need one more relationship. We must use another of the boundary conditions. The equations for E_x and E_y are no help, because all the E 's have only a z -component. So we must use the conditions on B . Let's try Eq. (33.29):

$$B_{x2} = B_{x1}.$$

From Eqs. (33.35) through (33.37),

$$B_{ri} = \frac{k_y E_i}{\omega}, \quad B_{ri} = \frac{k'_y E_r}{\omega'}, \quad B_{xt} = \frac{k''_y E_t}{\omega''}.$$

Recalling that $\omega'' = \omega' = \omega$ and $k''_y = k'_y = k_y$, we get that

$$E_0 + E'_0 = E''_0.$$

But this is just Eq. (33.48) all over again! We've just wasted time getting something we already knew.

We could try Eq. (33.30), $B_{z2} = B_{z1}$, but there are no z -components of \mathbf{B}' . So there's only one equation left: Eq. (33.31), $B_{y2} = B_{y1}$. For the three waves.

$$B_{yi} = -\frac{k_x E_i}{\omega}, \quad B_{yr} = -\frac{k'_x E_r}{\omega'}, \quad B_{yt} = -\frac{k''_x E_t}{\omega''}. \quad (33.49)$$

Putting for E_i , E_r , and E_t the wave expression for $x = 0$ (to be at the boundary), the boundary condition is

$$\frac{k_x}{\omega} E_0 e^{i(\omega t - k_y y)} + \frac{k'_x}{\omega'} E'_0 e^{i(\omega' t - k'_y y)} = \frac{k''_x}{\omega''} E''_0 e^{i(\omega'' t - k''_y y)}.$$

Again all ω 's and k_y 's are equal, so this reduces to

$$k_x E_0 + k'_x E'_0 = k''_x E''_0. \quad (33.50)$$

This gives us an equation for the E 's that is different from Eq. (33.48). With the two, we can solve for E'_0 and E''_0 . Remembering that $k'_x = -k_x$, we get

$$E'_0 = \frac{k_x - k''_x}{k_x + k''_x} E_0, \quad (33.51)$$

$$E''_0 = \frac{2k_x}{k_x + k''_x} E_0. \quad (33.52)$$

These, together with Eq. (33.45) or Eq. (33.46) for k''_x , give us what we wanted to know. We will discuss the consequences of this result in the next section.

If we begin with a wave polarized with its \mathbf{E} -vector *parallel* to the plane of incidence, \mathbf{E} will have both x - and y -components, as shown in Fig. 33-7. The algebra is straightforward but more complicated. (The work can be somewhat reduced by expressing things in this case in terms of the *magnetic* fields, which are all in the z -direction.) One finds that

$$|E'_0| = \frac{n_2^2 k_x - n_1^2 k''_x}{n_2^2 k_x + n_1^2 k''_x} |E_0| \quad (33.53)$$

and

$$|E''_0| = \frac{2n_1 n_2 k_x}{n_2^2 k_x + n_1^2 k''_x} |E_0|. \quad (33.54)$$

Let's see whether our results agree with those we got earlier. Equation (33.3) is the result we worked out in Chapter 35 of Volume I for the ratio of the intensity of the reflected wave to the intensity of the incident wave. Then, however, we were considering only *real* indexes. For real indexes (and k 's), we can write

$$k_x = k \cos \theta_i = \frac{\omega n_1}{c} \cos \theta_i,$$

$$k''_x = k'' \cos \theta_t = \frac{\omega n_2}{c} \cos \theta_t$$

Substituting in Eq. (33.51), we have

$$\frac{E'_0}{E_0} = \frac{n_1 \cos \theta_i - n_2 \cos \theta_t}{n_1 \cos \theta_i + n_2 \cos \theta_t}, \quad (33.55)$$

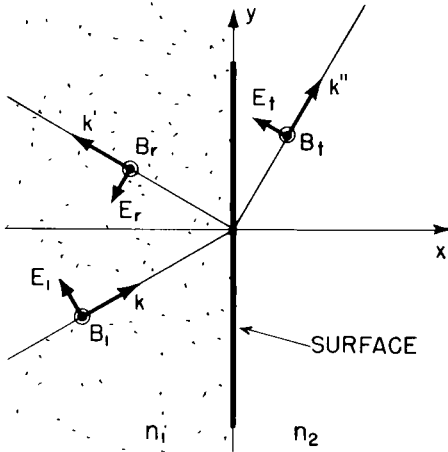


Fig. 33-7. Polarization of the waves when the \mathbf{E} -field of the incident wave is parallel to the plane of incidence.

which does not look the same as Eq. (33.3). It will, however, if we use Snell's law to get rid of the n 's. Setting $n_2 = n_1 \sin \theta_i / \sin \theta_t$, and multiplying the numerator and denominator by $\sin \theta_t$, we get

$$\frac{E'_0}{E_0} = \frac{\cos \theta_i \sin \theta_t - \sin \theta_i \cos \theta_t}{\cos \theta_i \sin \theta_t + \sin \theta_i \cos \theta_t}.$$

The numerator and denominator are just the sines of $(\theta_i - \theta_t)$ and $(\theta_i + \theta_t)$; we get

$$\frac{E'_0}{E_0} = \frac{\sin(\theta_i - \theta_t)}{\sin(\theta_i + \theta_t)}. \quad (33.56)$$

Since E'_0 and E_0 are in the same material, the intensities are proportional to the squares of the electric fields, and we get the same result as before. Similarly, Eq. (33.53) is the same as Eq. (33.4).

For waves which arrive at normal incidence, $\theta_i = 0$ and $\theta_t = 0$. Equation (33.56) gives $0/0$, which is not very useful. We can, however, go back to Eq. (33.55), which gives

$$\frac{I_r}{I_i} = \left(\frac{E'_0}{E_0} \right)^2 = \left(\frac{n_1 - n_2}{n_1 + n_2} \right)^2. \quad (33.57)$$

This result, naturally, applies for "either" polarization, since for normal incidence there is no special "plane of incidence."

33-5 Reflection from metals

We can now use our results to understand the interesting phenomenon of reflection from metals. Why is it that metals are shiny? We saw in the last chapter that metals have an index of refraction which, for some frequencies, has a large imaginary part. Let's see what we would get for the reflected intensity when light shines from air (with $n = 1$) onto a material with $n = -in_I$. Then Eq. (33.55) gives (for normal incidence)

$$\frac{E'_0}{E_0} = \frac{1 + in_I}{1 - in_I}.$$

For the *intensity* of the reflected wave, we want the square of the absolute values of E'_0 and E_0 :

$$\frac{I_r}{I_i} = \frac{|E'_0|^2}{|E_0|^2} = \frac{|1 + in_I|^2}{|1 - in_I|^2},$$

or

$$\frac{I_r}{I_i} = \frac{1 + n_I^2}{1 + n_I^2} = 1. \quad (33.58)$$

For a material with an index which is a pure imaginary number, there is 100 percent reflection!

Metals do not reflect 100 percent, but many do reflect visible light very well. In other words, the imaginary part of their indexes is very large. But we have seen that a large imaginary part of the index means a strong absorption. So there is a general rule that if *any* material gets to be a *very* good absorber at any frequency, the waves are strongly reflected at the surface and very little gets inside to be absorbed. You can see this effect with strong dyes. Pure crystals of the strongest dyes have a "metallic" shine. Probably you have noticed that at the edge of a bottle of purple ink the dried dye will give a golden metallic reflection, or that dried red ink will sometimes give a greenish metallic reflection. Red ink absorbs out the greens of *transmitted* light, so if the ink is very concentrated, it will exhibit a strong surface *reflection* for the frequencies of green light.

You can easily show this effect by coating a glass plate with red ink and letting it dry. If you direct a beam of white light at the back of the plate, as shown in Fig. 33-8, there will be a transmitted beam of red light and a reflected beam of green light.

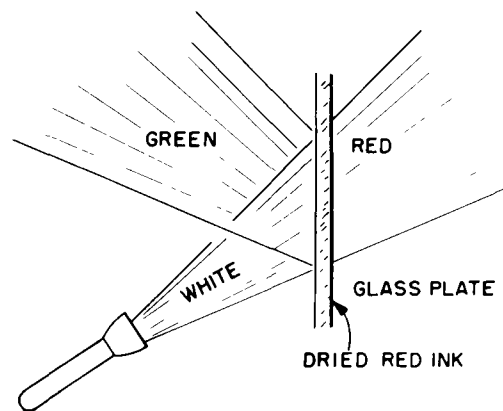


Fig. 33-8. A material which absorbs light strongly at the frequency ω also reflects light of that frequency.

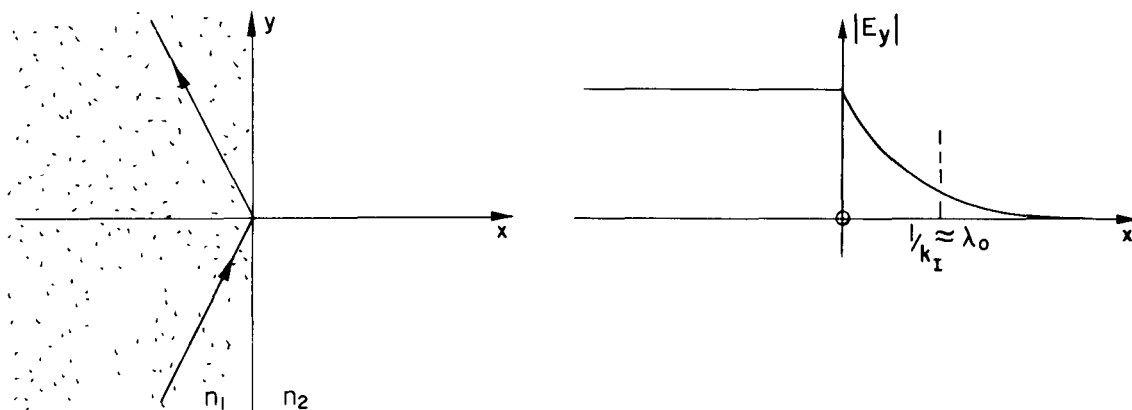


Fig. 33-9. Total internal reflection.

33-6 Total internal reflection

If light goes from a material like glass, with a real index n greater than 1, toward, say, air, with an index n_2 equal to 1, Snell's law says that

$$\sin \theta_t = n \sin \theta_i.$$

The angle θ_t of the transmitted wave becomes 90° when the incident angle θ_i is equal to the "critical angle" θ_c given by

$$n \sin \theta_c = 1. \quad (33.59)$$

What happens for θ_i greater than the critical angle? You know that there is total internal reflection. But how does that come about?

Let's go back to Eq (33.45) which gives the wave number k'' for the transmitted wave. We would have

$$k_x'^2 = \frac{k^2}{n^2} - k_y^2.$$

Now $k_y = k \sin \theta_i$ and $k = \omega n/c$, so

$$k_x'^2 = \frac{\omega^2}{c^2} (1 - n^2 \sin^2 \theta_i).$$

If $n \sin \theta_i$ is greater than one, $k_x'^2$ is *negative* and k_x' is a pure imaginary, say $\pm i k_I$. You know by now what that means! The "transmitted" wave (Eq. 33.34) will have the form

$$E_t = E_0'' e^{\pm k_I x} e^{i(\omega t - k_y y)}.$$

The wave amplitude either grows or drops off exponentially with increasing x . Clearly, what we want here is the negative sign. Then the *amplitude* of the wave to the right of the boundary will go as shown in Fig. 33-9. Notice that k_I is of the order ω/c —which is λ_0 , the free-space wavelength of the light. When light is totally reflected from the inside of a glass-air surface, there are fields in the air, but they extend beyond the surface only a distance of the order of the wavelength of the light

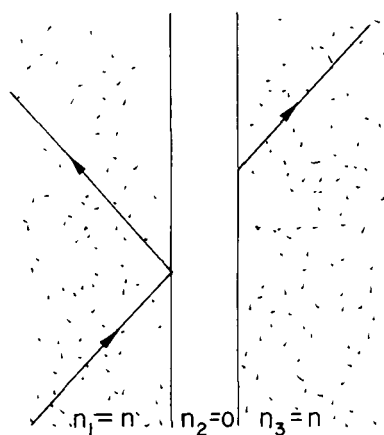


Fig. 33-10. If there is a small gap, internal reflection is not "total"; a transmitted wave appears beyond the gap.

We can now see how to answer the following question: If a light wave in glass arrives at the surface at a large enough angle, it is reflected, if another piece of glass is brought up to the surface (so that the "surface" in effect disappears) the light is transmitted. Exactly when does this happen? Surely there must be continuous change from total reflection to no reflection! The answer, of course, is that if the air gap is so small that the exponential tail of the wave in the air has an appreciable strength at the second piece of glass, it will shake the electrons there and generate a new wave, as shown in Fig. 33-10. Some light will be transmitted. (Clearly, our solution is incomplete, we should solve all the equations again for a thin layer of air between two regions of glass.)

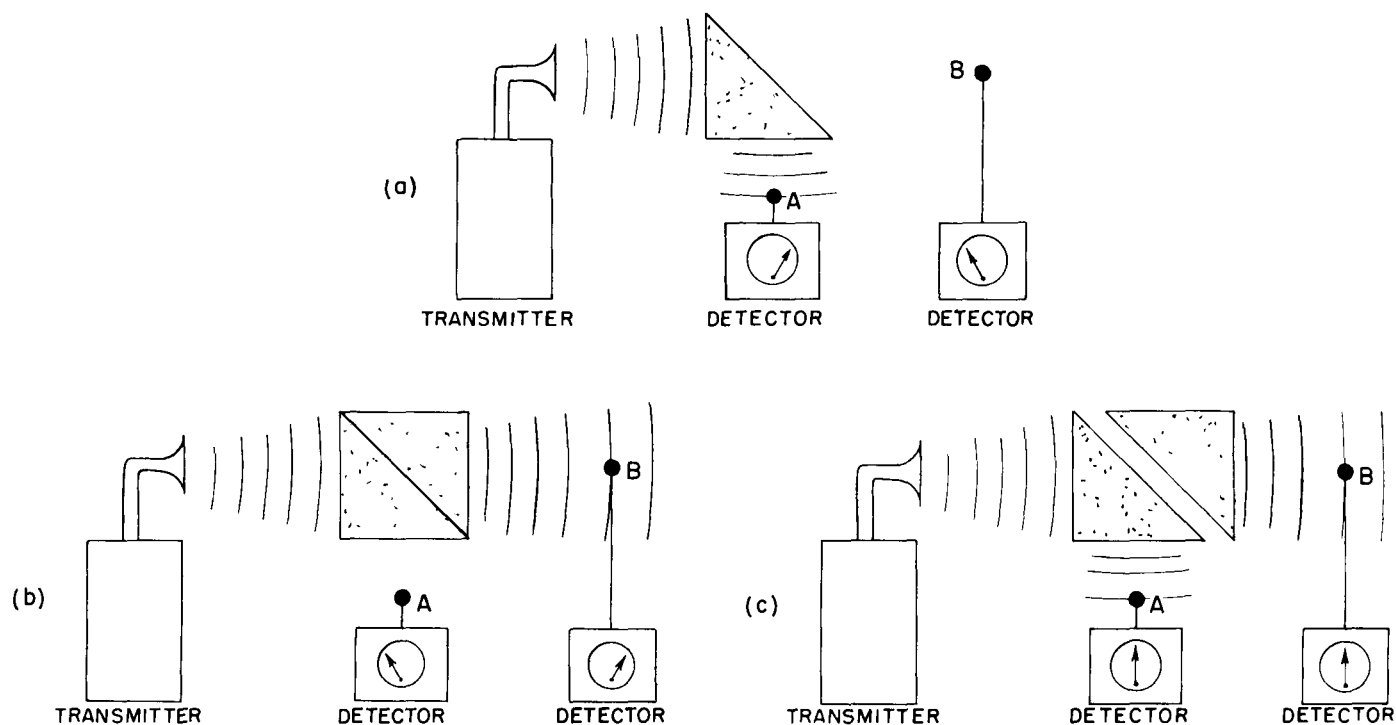


Fig. 33-11. A demonstration of the penetration of internally reflected waves.

This transmission effect can be observed with ordinary light only if the air gap is very small (of the order of the wavelength of light, like 10^{-5} cm), but it is easily demonstrated with three-centimeter waves. Then the exponentially decreasing field extends several centimeters. A microwave apparatus that shows the effect is drawn in Fig. 33-11. Waves from a small three-centimeter transmitter are directed at a 45° prism of paraffin. The index of refraction of paraffin for these frequencies is 1.50, and therefore the critical angle is 41.5° . So the wave is totally reflected from the 45° face and is picked up by detector *A*, as indicated in Fig. 33-11(a). If a second paraffin prism is placed in contact with the first, as shown in part (b) of the figure, the wave passes straight through and is picked up at detector *B*. If a gap of a few centimeters is left between the two prisms, as in part (c), there are both transmitted and reflected waves. The electric field outside the 45° face of the prism in Fig. 33-11(a) can also be shown by bringing detector *B* to within a few centimeters of the surface.

The Magnetism of Matter

34-1 Diamagnetism and paramagnetism

In this chapter we are going to talk about the magnetic properties of materials. The material which has the most striking magnetic properties is, of course, iron. Similar magnetic properties are shared also by the elements nickel, cobalt, and—at sufficiently low temperatures (below 16°C)—by gadolinium, as well as by a number of peculiar alloys. That kind of magnetism, called *ferromagnetism*, is sufficiently striking and complicated that we will discuss it in a special chapter. However, all ordinary substances do show some magnetic effects, although very small ones—a thousand to a million times less than the effects in ferromagnetic materials. Here we are going to describe ordinary magnetism, that is to say, the magnetism of substances other than the ferromagnetic ones.

This small magnetism is of two kinds. Some materials are *attracted* toward magnetic fields; others are *repelled*. Unlike the electrical effect in matter, which always causes dielectrics to be attracted, there are two signs to the magnetic effect. These two signs can be easily shown with the help of a strong electromagnet which has one sharply pointed pole piece and one flat pole piece, as drawn in Fig. 34-1. The magnetic field is much stronger near the pointed pole than near the flat pole. If a small piece of material is fastened to a long string and suspended between the poles, there will, in general, be a small force on it. This small force can be seen by the slight displacement of the hanging material when the magnet is turned on. The few ferromagnetic materials are attracted very strongly toward the pointed pole; all other materials feel only a very weak force. Some are weakly attracted to the pointed pole; and some are weakly repelled.

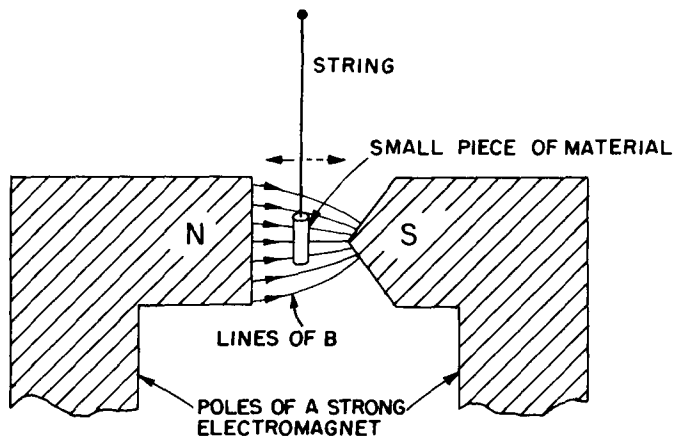


Fig. 34-1. A small cylinder of bismuth is weakly repelled by the sharp pole; a piece of aluminum is attracted.

The effect is most easily seen with a small cylinder of bismuth, which is *repelled* from the high-field region. Substances which are repelled in this way are called *diamagnetic*. Bismuth is one of the strongest diamagnetic materials, but even with it, the effect is still quite weak. Diamagnetism is always very weak. If a small piece of aluminum is suspended between the poles, there is also a weak force, but *toward* the pointed pole. Substances like aluminum are called *paramagnetic*. (In such an experiment, eddy-current forces arise when the magnet is turned on and off, and these can give off strong impulses. You must be careful to look for the net displacement after the hanging object settles down.)

34-1 Diamagnetism and paramagnetism

34-2 Magnetic moments and angular momentum

34-3 The precession of atomic magnets

34-4 Diamagnetism

34-5 Larmor's theorem

34-6 Classical physics gives neither diamagnetism nor paramagnetism

34-7 Angular momentum in quantum mechanics

34-8 The magnetic energy of atoms

Review: Section 15-1, "The forces on a current loop; energy of a dipole."

We want now to describe briefly the mechanisms of these two effects. First, in many substances the atoms have no permanent magnetic moments, or rather, all the magnets within each atom balance out so that the *net* moment of the atom is zero. The electron spins and orbital motions all exactly balance out, so that any particular atom has no average magnetic moment. In these circumstances, when you turn on a magnetic field little extra currents are generated inside the atom by induction. According to Lenz's law, these currents are in such a direction as to oppose the increasing field. So the induced magnetic moments of the atoms are directed *opposite* to the magnetic field. This is the mechanism of diamagnetism.

Then there are some substances for which the atoms do have a permanent magnetic moment—in which the electron spins and orbits have a net circulating current that is not zero. So besides the diamagnetic effect (which is always present), there is also the possibility of lining up the individual atomic magnetic moments. In this case, the moments try to line up *with* the magnetic field (in the way the permanent dipoles of a dielectric are lined up by the electric field), and the induced magnetism tends to enhance the magnetic field. These are the paramagnetic substances. Paramagnetism is generally fairly weak because the lining-up forces are relatively small compared with the forces from the thermal motions which try to derange the order. It also follows that paramagnetism is usually sensitive to the temperature. (The paramagnetism arising from the spins of the electrons responsible for conduction in a metal constitutes an exception. We will not be discussing this phenomenon here.) For ordinary paramagnetism, the lower the temperature, the stronger the effect. There is more lining-up at low temperatures when the deranging effects of the collisions are less. Diamagnetism, on the other hand, is more or less independent of the temperature. In any substance with built-in magnetic moments there is a diamagnetic as well as a paramagnetic effect, but the paramagnetic effect usually dominates.

In Chapter 11 we described a *ferroelectric* material, in which all the electric dipoles get lined up by their own mutual electric fields. It is also possible to imagine the magnetic analog of ferroelectricity, in which all the atomic moments would line up and lock together. If you make calculations of how this should happen, you will find that because the magnetic forces are so much smaller than the electric forces, thermal motions should knock out this alignment even at temperatures as low as a few tenths of a degree Kelvin. So it would be impossible at room temperature to have any permanent lining up of the magnets.

On the other hand, this is exactly what does happen in iron—it does get lined up. There is an effective force between the magnetic moments of the different atoms of iron which is much, much greater than the *direct magnetic* interaction. It is an indirect effect which can be explained only by quantum mechanics. It is about ten thousand times stronger than the direct magnetic interaction, and is what lines up the moments in ferromagnetic materials. We discuss this special interaction in a later chapter.

Now that we have tried to give you a qualitative explanation of diamagnetism and paramagnetism, we must correct ourselves and say that *it is not possible* to understand the magnetic effects of materials in any honest way from the point of view of classical physics. Such magnetic effects are a *completely quantum-mechanical phenomenon*. It is, however, possible to make some phoney classical arguments and to get some idea of what is going on. We might put it this way. You can make some classical arguments and get guesses as to the behavior of the material, but these arguments are not “legal” in any sense because it is absolutely essential that quantum mechanics be involved in every one of these magnetic phenomena. On the other hand, there are situations, such as in a plasma or a region of space with many free electrons, where the electrons do obey the laws of classical mechanics. And in those circumstances, some of the theorems from classical magnetism are worth while. Also, the classical arguments are of some value for historical reasons. The first few times that people were able to guess at the meaning and behavior of magnetic materials, they used classical arguments. Finally, as we have already illustrated, classical mechanics can give us some useful guesses

as to what might happen—even though the really honest way to study this subject would be to learn quantum mechanics first and then to understand the magnetism in terms of quantum mechanics.

On the other hand, we don't want to wait until we learn quantum mechanics inside out to understand a simple thing like diamagnetism. We will have to lean on the classical mechanics as kind of half showing what happens, realizing, however, that the arguments are really not correct. We therefore make a series of theorems about classical magnetism that will confuse you because they will prove different things. Except for the last theorem, every one of them will be wrong. Furthermore, they will all be wrong as a description of the physical world, because quantum mechanics is left out.

34-2 Magnetic moments and angular momentum

The first theorem we want to prove from classical mechanics is the following: If an electron is moving in a circular orbit (for example, revolving around a nucleus under the influence of a central force), there is a definite ratio between the magnetic moment and the angular momentum. Let's call \mathbf{J} the angular momentum and $\boldsymbol{\mu}$ the magnetic moment of the electron in the orbit. The magnitude of the angular momentum is the mass of the electron times the velocity times the radius (See Fig. 34-2.) It is directed perpendicular to the plane of the orbit.

$$\mathbf{J} = m\mathbf{v}r. \quad (34.1)$$

(This is, of course, a nonrelativistic formula, but it is a good approximation for atoms, because for the electrons involved v/c is generally of the order of $e^2/\hbar c = 1/137$, or about 1 percent.)

The magnetic moment of the same orbit is the current times the area. (See Section 14-5.) The current is the charge per unit time which passes any point on the orbit, namely, the charge q times the frequency of rotation. The frequency is the velocity divided by the circumference of the orbit; so

$$I = q \frac{v}{2\pi r}.$$

The area is πr^2 , so the magnetic moment is

$$\boldsymbol{\mu} = \frac{qvr}{2}. \quad (34.2)$$

It is also directed perpendicular to the plane of the orbit. So \mathbf{J} and $\boldsymbol{\mu}$ are in the same direction:

$$\boldsymbol{\mu} = \frac{q}{2m} \mathbf{J} \text{ (orbit)}. \quad (34.3)$$

Their ratio depends neither on the velocity nor on the radius. For any particle moving in a circular orbit the magnetic moment is equal to $q/2m$ times the angular momentum. For an electron, the charge is negative—we can call it $-q_e$; so for an electron

$$\boldsymbol{\mu} = -\frac{q_e}{2m} \mathbf{J} \text{ (electron orbit)}. \quad (34.4)$$

That's what we would expect classically and, miraculously enough, it is also true quantum-mechanically. It's one of those things. However, if you keep going with the classical physics, you find other places where it gives the wrong answers, and it is a great game to try to remember which things are right and which things are wrong. We might as well give you immediately what is true *in general* in quantum mechanics. First, Eq. (34.4) is true for *orbital motion*, but that's not the only magnetism that exists. The electron also has a spin rotation about its own axis (something like the earth rotating on its axis), and as a result of that spin it has both an angular momentum and a magnetic moment. But for reasons that are purely quantum-mechanical—there is no classical explanation—the ratio of $\boldsymbol{\mu}$

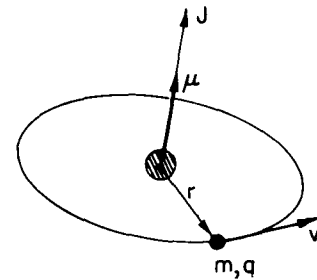


Fig. 34-2. For any circular orbit the magnetic moment $\boldsymbol{\mu}$ is $q/2m$ times the angular momentum \mathbf{J} .

to \mathbf{J} for the electron spin is twice as large as it is for orbital motion of the spinning electron:

$$\boldsymbol{\mu} = -\frac{q_e}{m} \mathbf{J} \text{ (electron spin).} \quad (34.5)$$

In any atom there are, generally speaking, several electrons and some combination of spin and orbit rotations which builds up a total angular momentum and a total magnetic moment. Although there is no classical reason why it should be so, it is *always true* in quantum mechanics that (for an isolated atom) the direction of the magnetic moment is exactly opposite to the direction of the angular momentum. The ratio of the two is not necessarily either $-q_e/m$ or $-q_e/2m$, but somewhere in between, because there is a mixture of the contributions from the orbits and the spins. We can write

$$\boldsymbol{\mu} = -g \left(\frac{q_e}{2m} \right) \mathbf{J}, \quad (34.6)$$

where g is a factor which is characteristic of the state of the atom. It would be 1 for a pure orbital moment, or 2 for a pure spin moment, or some other number in between for a complicated system like an atom. This formula does not, of course, tell us very much. It says that the magnetic moment is *parallel* to the angular momentum, but can have any magnitude. The form of Eq. (34.6) is convenient, however, because g —called the “Landé g -factor”—is a dimensionless constant whose magnitude is of the order of one. It is one of the jobs of quantum mechanics to predict the g -factor for any particular atomic state.

You might also be interested in what happens in nuclei. In nuclei there are protons and neutrons which may move around in some kind of orbit and at the same time, like an electron, have an intrinsic spin. Again the magnetic moment is parallel to the angular momentum. Only now the order of magnitude of the ratio of the two is what you would expect for a *proton* going around in a circle, with m in Eq. (34.3) equal to the *proton* mass. Therefore it is usual to write for nuclei

$$\boldsymbol{\mu} = g \left(\frac{q_e}{2m_p} \right) \mathbf{J}, \quad (34.7)$$

where m_p is the mass of the proton, and g —called the *nuclear* g -factor—is a number near one, to be determined for each nucleus.

Another important difference for a nucleus is that the *spin* magnetic moment of the proton does *not* have a g -factor of 2, as the electron does. For a proton, $g = 2(2.79)$. Surprisingly enough, the *neutron* also has a spin magnetic moment, and its magnetic moment relative to its angular momentum is $2(-1.93)$. The neutron, in other words, is not exactly “neutral” in the magnetic sense. It is like a little magnet, and it has the kind of magnetic moment that a rotating *negative* charge would have.

34-3 The precession of atomic magnets

One of the consequences of having the magnetic moment proportional to the angular momentum is that an atomic magnet placed in a magnetic field will *precess*. First we will argue classically. Suppose that we have the magnetic moment $\boldsymbol{\mu}$ suspended freely in a uniform magnetic field. It will feel a torque $\boldsymbol{\tau}$, equal to $\boldsymbol{\mu} \times \mathbf{B}$, which tries to bring it in line with the field direction. But the atomic magnet is a gyroscope—it has the angular momentum \mathbf{J} . Therefore the torque due to the magnetic field will not cause the magnet to line up. Instead, the magnet will *precess*, as we saw when we analyzed a gyroscope in Chapter 20 of Volume I. The angular momentum—and with it the magnetic moment—precesses about an axis parallel to the magnetic field. We can find the rate of precession by the same method we used in Chapter 20 of the first volume.

Suppose that in a small time Δt the angular momentum changes from \mathbf{J} to \mathbf{J}' , as drawn in Fig. 34-3, staying always at the same angle θ with respect to the direction of the magnetic field \mathbf{B} . Let's call ω_p the angular velocity of the precession, so that in the time Δt the angle of *precession* is $\omega_p \Delta t$. From the geometry of the

figure, we see that the change of angular momentum in the time Δt is

$$\Delta J = (J \sin \theta)(\omega_p \Delta t).$$

So the rate of change of the angular momentum is

$$\frac{dJ}{dt} = \omega_p J \sin \theta, \quad (34.8)$$

which must be equal to the torque:

$$\tau = \mu B \sin \theta. \quad (34.9)$$

The angular velocity of precession is then

$$\omega_p = \frac{\mu}{J} B. \quad (34.10)$$

Substituting μ/J from Eq. (34.6), we see that for an atomic system

$$\omega_p = g \frac{q_e B}{2m}, \quad (34.11)$$

the precession frequency is proportional to B . It is handy to remember that for an atom (or electron)

$$f_p = \frac{\omega_p}{2\pi} = (1.4 \text{ megacycles/gauss})gB, \quad (34.12)$$

and that for a nucleus

$$f_p = \frac{\omega_p}{2\pi} = (0.76 \text{ kilocycles/gauss})gB. \quad (34.13)$$

(The formulas for atoms and nuclei are different only because of the different conventions for g for the two cases.)

According to the *classical* theory, then, the electron orbits—and spins—in an atom should precess in a magnetic field. Is it also true quantum-mechanically? It is essentially true, but the meaning of the “precession” is different. In quantum mechanics one cannot talk about the *direction* of the angular momentum in the same sense as one does classically, nevertheless, there is a very close analogy—so close that we continue to call it “precession.” We will discuss it later when we talk about the quantum-mechanical point of view.

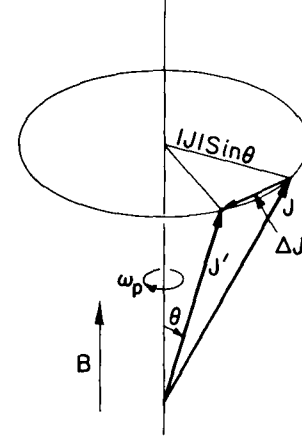


Fig. 34-3. An object with angular momentum J and a parallel magnetic moment μ placed in a magnetic field B precesses with the angular velocity ω_p .

34-4 Diamagnetism

Next we want to look at *diamagnetism* from the classical point of view. It can be worked out in several ways, but one of the nice ways is the following. Suppose that we slowly turn on a magnetic field in the vicinity of an atom. As the magnetic field changes an *electric* field is generated by magnetic induction. From Faraday's law, the line integral of E around any closed path is the rate of change of the magnetic flux through the path. Suppose we pick a path Γ which is a circle of radius r concentric with the center of the atom, as shown in Fig. 34-4. The average tangential electric field E around this path is given by

$$E 2\pi r = - \frac{d}{dt} (B \pi r^2),$$

and there is a circulating electric field whose strength is

$$E = - \frac{r}{2} \frac{dB}{dt}.$$

The induced electric field acting on an electron in the atom produces a torque equal to $-q_e E r$, which must equal the rate of change of the angular momentum dJ/dt :

$$\frac{dJ}{dt} = \frac{q_e r^2}{2} \frac{dB}{dt}. \quad (34.14)$$

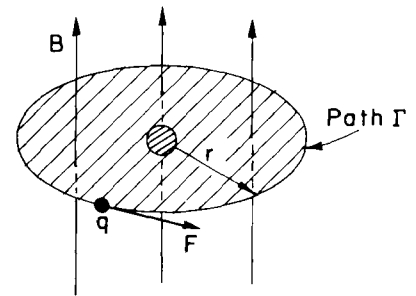


Fig. 34-4. The induced electric forces on the electrons in an atom.

Integrating with respect to time from zero field, we find that the change in angular momentum due to turning on the field is

$$\Delta J = \frac{q_e r^2}{2} B. \quad (34.15)$$

This is the extra angular momentum from the twist given to the electrons as the field is turned on.

This added angular momentum makes an extra magnetic moment which, because it is an *orbital* motion, is just $-q_e/2m$ times the angular momentum. The induced diamagnetic moment is

$$\Delta\mu = -\frac{q_e}{2m} \Delta J = -\frac{q_e^2 r^2}{4m} B. \quad (34.16)$$

The minus sign (as you can see is right by using Lenz's law) means that the added moment is opposite to the magnetic field.

We would like to write Eq. (34.16) a little differently. The r^2 which appears is the radius from an axis through the atom parallel to B , so if B is along the z -direction, it is $x^2 + y^2$. If we consider spherically symmetric atoms (or average over atoms with their natural axes in all directions) the average of $x^2 + y^2$ is $2/3$ of the average of the square of the true radial distance from the center *point* of the atom. It is therefore usually more convenient to write Eq. (34.16) as

$$\Delta\mu = -\frac{q_e^2}{6m} \langle r^2 \rangle_{av} B. \quad (34.17)$$

In any case, we have found an induced atomic moment proportional to the magnetic field B and opposing it. This is diamagnetism of matter. It is this magnetic effect that is responsible for the small force on a piece of bismuth in a nonuniform magnetic field. (You could compute the force by working out the energy of the induced moments in the field and seeing how the energy changes as the material is moved into or out of the high-field region.)

We are still left with the problem: What is the mean square radius, $\langle r^2 \rangle_{av}$? Classical mechanics cannot supply an answer. We must go back and start over with quantum mechanics. In an atom we cannot really say where an electron is, but only know the probability that it will be at some place. If we interpret $\langle r^2 \rangle_{av}$ to mean the average of the square of the distance from the center for the probability distribution, the diamagnetic moment given by quantum mechanics is just the same as formula (34.17). This equation, of course, is the moment for one electron. The total moment is given by the sum over all the electrons in the atom. The surprising thing is that the classical argument and quantum mechanics give the same answer, although, as we shall see, the classical argument that gives Eq. (34.17) is not really valid in classical mechanics.

The same diamagnetic effect occurs even when an atom already has a permanent moment. Then the system will precess in the magnetic field. As the whole atom precesses, it takes up an additional small angular velocity, and that slow turning gives a small current which represents a correction to the magnetic moment. This is just the diamagnetic effect represented in another way. But we don't really have to worry about that when we talk about paramagnetism. If the diamagnetic effect is first computed, as we have done here, we don't have to worry about the fact that there is an extra little current from the precession. That has already been included in the diamagnetic term.

34-5 Larmor's theorem

We can already conclude something from our results so far. First of all, in the classical theory the moment μ was always proportional to J , with a given constant of proportionality for a particular atom. There wasn't any spin of the electrons, and the constant of proportionality was always $-q_e/2m$; that is to say, in Eq. (34.6) we should set $g = 1$. The ratio of μ to J was independent of the internal motion of the electrons. Thus, according to the classical theory, all systems

of electrons would precess with *the same* angular velocity. (This is *not* true in quantum mechanics.) This result is related to a theorem in classical mechanics that we would now like to prove. Suppose we have a group of electrons which are all held together by attraction toward a central point—as the electrons are attracted by a nucleus. The electrons will also be interacting with each other, and can, in general, have complicated motions. Suppose you have solved for the motions with *no* magnetic field and then want to know what the motions would be *with* a weak magnetic field. The theorem says that the motion with a weak magnetic field is always one of the no-field solutions with an added rotation, about the axis of the field, with the angular velocity $\omega_L = q_e B / 2m$. (This is the same as ω_p , if $g = 1$.) There are, of course, many possible motions. The point is that for every motion without the magnetic field there is a corresponding motion in the field, which is the original motion plus a uniform rotation. This is called Larmor's theorem, and ω_L is called the *Larmor frequency*.

We would like to show how the theorem can be proved, but we will let you work out the details. Take, first, one electron in a central force field. The force on it is just $F(r)$, directed toward the center. If we now turn on a uniform magnetic field, there is an additional force, $qv \times B$; so the total force is

$$F(r) + qv \times B. \quad (34.18)$$

Now let's look at the same system from a coordinate system rotating with angular velocity ω about an axis through the center of force and parallel to B . This is no longer an inertial system, so we have to put in the proper pseudoforces—the centrifugal and Coriolis forces we talked about in Chapter 19 of Volume I. We found there that in a frame rotating with angular velocity ω , there is an apparent *tangential* force proportional to v_r , the radial component of velocity:

$$F_t = -2m\omega v_r, \quad (34.19)$$

And there is an apparent radial force which is given by

$$F_r = m\omega^2 r + 2m\omega v_t, \quad (34.20)$$

where v_t is the tangential component of the velocity, measured *in* the rotating frame. (The radial component v_r for rotating and inertial frames is the same.)

Now for small enough angular velocities (that is, if $\omega r \ll v_t$), we can neglect the first term (centrifugal) in Eq. (34.20) in comparison with the second (Coriolis). Then Eqs. (34.19) and (34.20) can be written together as

$$F = -2m\omega \times v. \quad (34.21)$$

If we now *combine* a rotation and a magnetic field, we must add the force in Eq. (34.21) to that in Eq. (34.18). The total force is

$$F(r) + qv \times B + 2mv \times \omega \quad (34.22)$$

[we reverse the cross product and the sign of Eq. (34.21) to get the last term]. Looking at our result, we see that if

$$2m\omega = -qB$$

the two terms on the right cancel, and in the moving frame the only force is $F(r)$. The motion of the electron is just the same as with no magnetic field—and, of course, no rotation. We have proved Larmor's theorem for one electron. Since the proof assumes a small ω , it also means that the theorem is true only for weak magnetic fields. The only thing we could ask you to improve on is to take the case of many electrons mutually interacting with each other, but all in the same central field, and prove the same theorem. So no matter how complex an atom is, if it has a central field the theorem is true. But that's the end of the classical mechanics, because it isn't true in fact that the motions precess in that way. The precession frequency ω_p of Eq. (34.11) is only equal to ω_L if g happens to be equal to 1.

34-6 Classical physics gives neither diamagnetism nor paramagnetism

Now we would like to demonstrate that according to classical mechanics there can be no diamagnetism and no paramagnetism at all. It sounds crazy—first, we have proved that there are paramagnetism, diamagnetism, precessing orbits, and so on, and now we are going to prove that it is all wrong. Yes!—We are going to prove that *if you follow the classical mechanics far enough*, there are no such magnetic effects—they all cancel out. If you start a classical argument in a certain place and don't go far enough, you can get any answer you want. But the only legitimate and correct proof shows that there is no magnetic effect whatever.

It is a consequence of classical mechanics that if you have any kind of system—a gas with electrons, protons, and whatever—kept in a box so that the whole thing can't turn, there will be no magnetic effect. It is possible to have a magnetic effect if you have an isolated system, like a star held together by itself, which can start rotating when you put on the magnetic field. But if you have a piece of material that is held in place so that it can't start spinning, then there will be no magnetic effects. What we mean by holding down the spin is summarized this way: At a given temperature we suppose that there is *only one state* of thermal equilibrium. The theorem then says that if you turn on a magnetic field and wait for the system to get into thermal equilibrium, there will be no paramagnetism or diamagnetism—there will be no induced magnetic moment. Proof: According to statistical mechanics, the probability that a system will have any given state of motion is proportional to $e^{-U/kT}$, where U is the energy of that motion. Now what is the energy of motion? For a particle moving in a constant magnetic field, the energy is the ordinary potential energy plus $mv^2/2$, with nothing additional for the magnetic field. [You know that the forces from electromagnetic fields are $q(\mathbf{E} + \mathbf{v} \times \mathbf{B})$, and that the rate of work $\mathbf{F} \cdot \mathbf{v}$ is just $q\mathbf{E} \cdot \mathbf{v}$, which is not affected by the magnetic field.] So the energy of a system, whether it is in a magnetic field or not, is always given by the kinetic energy plus the potential energy. Since the probability of any motion depends only on the energy—that is, on the velocity and position—it is the same whether or not there is a magnetic field. For *thermal* equilibrium, therefore, the magnetic field has no effect. If we have one system in a box, and then have another system in a second box, this time with a magnetic field, the probability of any particular velocity at any point in the first box is the same as in the second. If the first box has no average circulating current (which it will not have if it is in equilibrium with the stationary walls), there is no average magnetic moment. Since in the second box all the motions are the same, there is no average magnetic moment there either. Hence, if the temperature is kept constant and thermal equilibrium is re-established after the field is turned on, there can be no magnetic moment induced by the field—according to classical mechanics. We can only get a satisfactory understanding of magnetic phenomena from quantum mechanics.

Unfortunately, we cannot assume that you have a thorough understanding of quantum mechanics, so this is hardly the place to discuss the matter. On the other hand, we don't always have to learn something first by learning the exact rules and then by learning how they are applied in different cases. Almost every subject that we have taken up in this course has been treated in a different way. In the case of electricity, we wrote the Maxwell equations on "Page One" and then deduced all the consequences. That's one way. But we will *not* now try to begin a new "Page One," writing the equations of quantum mechanics and deducing everything from them. We will just have to tell you some of the consequences of quantum mechanics, before you learn where they come from. So here we go.

34-7 Angular momentum in quantum mechanics

We have already given you a relation between the magnetic moment and the angular momentum. That's pleasant. But what do the magnetic moment and the angular momentum *mean* in quantum mechanics? In quantum mechanics it turns out to be best to define things like magnetic moments in terms of the other concepts such as energy, in order to make sure that one knows what it means. Now,

it is easy to define a magnetic moment in terms of energy, because the energy of a moment in a magnetic field is, in the classical theory, $\boldsymbol{\mu} \cdot \mathbf{B}$. Therefore, the following definition has been taken in quantum mechanics: If we calculate the energy of a system in a magnetic field and we find that it is proportional to the field strength (for small field), the coefficient is called the component of magnetic moment in the direction of the field. (We don't have to get so elegant for our work now, we can still think of the magnetic moment in the ordinary, to some extent classical, sense.)

Now we would like to discuss the idea of angular momentum in quantum mechanics—or rather, the characteristics of what, in quantum mechanics, is called angular momentum. You see, when you go to new kinds of laws, you can't just assume that each word is going to mean exactly the same thing. You may think, say, “Oh, I know what angular momentum is. It's that thing that is changed by a torque.” But what's a torque? In quantum mechanics we have to have new definitions of old quantities. It would, therefore, be legally best to call it by some other name such as “quantangular momentum,” or something like that, because it is the angular momentum as defined in quantum mechanics. But if we can find a quantity in quantum mechanics which is identical to our old idea of angular momentum when the system becomes large enough, there is no use in inventing an extra word. We might as well just call it angular momentum. With that understanding, this odd thing that we are about to describe *is* angular momentum. It is the thing which in a large system we recognize as angular momentum in classical mechanics.

First, we take a system in which angular momentum is conserved, such as an atom all by itself in empty space. Now such a thing (like the earth spinning on its axis) could, in the ordinary sense, be spinning around any axis one wished to choose. And for a given spin, there could be many different “states,” all of the same energy, each “state” corresponding to a particular direction of the axis of the angular momentum. So in the classical theory, with a given angular momentum, there is an infinite number of possible states, all of the same energy.

It turns out in quantum mechanics, however, that several strange things happen. First, the number of states in which such a system *can exist* is limited—there is only a finite number. If the system is small, the finite number is very small, and if the system is large, the finite number gets very, very large. Second, we *cannot* describe a “state” by giving the *direction* of its angular momentum, but only by giving the *component* of the angular momentum along some direction—say in the *z*-direction. Classically, an object with a given total angular momentum J could have, for its *z*-component, any value from $+J$ to $-J$. But quantum-mechanically, the *z*-component of angular momentum can have only certain discrete values. Any given system—a particular atom, or a nucleus, or anything—with a given energy, has a characteristic number j , and its *z*-component of angular momentum can only be one of the following set of values:

$$\begin{array}{l} j\hbar \\ (j-1)\hbar \\ (j-2)\hbar \\ \vdots \\ -(j-2)\hbar \\ -(j-1)\hbar \\ -j\hbar \end{array} \quad (34.23)$$

The largest *z*-component is j times \hbar ; the next smaller is one unit of \hbar less, and so on down to $-j\hbar$. The number j is called “the spin of the system.” (Some people call it the “total angular momentum quantum number”; but we'll call it the “spin.”)

You may be worried that what we are saying can only be true for some “special” *z*-axis. But that is not so. For a system whose spin is j , the component of angular momentum along *any* axis can have only one of the values in (34.23). Although it is quite mysterious, we ask you just to accept it for the moment. We

will come back and discuss the point later. You may at least be pleased to hear that the z -component goes from some number to minus the *same* number, so that we at least don't have to decide which is the plus direction of the z -axis. (Certainly, if we said that it went from $+j$ to minus a different amount, that would be infinitely mysterious, because we wouldn't have been able to define the z -axis, pointing the other way.)

Now if the z -component of angular momentum must go down by integers from $+j$ to $-j$, then j must be an integer. No! Not quite; twice j must be an integer. It is only the *difference* between $+j$ and $-j$ that must be an integer. So, in general, the spin j is either an integer or a half-integer, depending on whether $2j$ is even or odd. Take, for instance, a nucleus like lithium, which has a spin of three-halves, $j = 3/2$. Then the angular momentum around the z -axis, in units of \hbar , is one of the following:

$$\begin{aligned} &+3/2 \\ &+1/2 \\ &-1/2 \\ &-3/2. \end{aligned}$$

There are four possible states, each of the same energy, if the nucleus is in empty space with no external fields. If we have a system whose spin is two, then the z -component of angular momentum has only the values, in units of \hbar ,

$$\begin{aligned} &2 \\ &1 \\ &0 \\ &-1 \\ &-2. \end{aligned}$$

If you count how many states there are for a given j , there are $(2j + 1)$ possibilities. In other words, if you tell me the energy and also the spin j , it turns out that there are exactly $(2j + 1)$ states with that energy, each state corresponding to one of the different possible values of the z -component of the angular momentum.

We would like to add one other fact. If you pick out any atom of known j at random and measure the z -component of the angular momentum, then you may get any one of the possible values, and each of the values is *equally* likely. All of the states are in fact single states, and each is just as good as any other. Each one has the same "weight" in the world. (We are assuming that nothing has been done to sort out a special sample.) This fact has, incidentally, a simple classical analog. If you ask the same question classically: What is the likelihood of a particular z -component of angular momentum if you take a random sample of systems, all with the same total angular momentum?—the answer is that all values from the maximum to the minimum are equally likely. (You can easily work that out.) The classical result corresponds to the equal probability of the $(2j + 1)$ possibilities in quantum mechanics.

From what we have so far, we can get another interesting and somewhat surprising conclusion. In certain classical calculations the quantity that appears in the final result is the *square* of the magnitude of the angular momentum \mathbf{J} —in other words, $\mathbf{J} \cdot \mathbf{J}$. It turns out that it is often possible to *guess* at the correct quantum-mechanical formula by using the classical calculation and the following simple rule: Replace $J^2 = \mathbf{J} \cdot \mathbf{J}$ by $j(j + 1)\hbar^2$. This rule is commonly used, and usually gives the correct result, but *not* always. We can give the following argument to show why you might expect this rule to work.

The scalar product $\mathbf{J} \cdot \mathbf{J}$ can be written as

$$\mathbf{J} \cdot \mathbf{J} = J_x^2 + J_y^2 + J_z^2.$$

Since it is a scalar, it should be the same for any orientation of the spin. Suppose we pick samples of any given atomic system at random and make measurements of J_x^2 , or J_y^2 , or J_z^2 , the *average value* should be the same for each. (There is no special distinction for any one of the directions.) Therefore, the average of $\mathbf{J} \cdot \mathbf{J}$ is just

equal to three times the average of any component squared, say of J_z^2 ;

$$\langle \mathbf{J} \cdot \mathbf{J} \rangle_{\text{av}} = 3 \langle J_z^2 \rangle.$$

But since $\mathbf{J} \cdot \mathbf{J}$ is the same for all orientations, its average is, of course, just its constant value; we have

$$\mathbf{J} \cdot \mathbf{J} = 3 \langle J_z^2 \rangle_{\text{av}}. \quad (34.24)$$

If we now say that we will use the same equation for quantum mechanics, we can easily find $\langle J_z^2 \rangle_{\text{av}}$. We just have to take the sum of the $(2j + 1)$ possible values of J_z^2 , and divide by the total number;

$$\langle J_z^2 \rangle_{\text{av}} = \frac{j^2 + (j-1)^2 + \cdots + (-j+1)^2 + (-j)^2}{2j+1} \hbar^2. \quad (34.25)$$

For a system with a spin of $3/2$, it goes like this:

$$\langle J_z^2 \rangle_{\text{av}} = \frac{(3/2)^2 + (1/2)^2 + (-1/2)^2 + (-3/2)^2}{4} \hbar^2 = \frac{5}{4} \hbar^2.$$

We conclude that

$$\mathbf{J} \cdot \mathbf{J} = 3 \langle J_z^2 \rangle_{\text{av}} = 3 \frac{5}{4} \hbar^2 = \frac{3}{2} \left(\frac{3}{2} + 1 \right) \hbar^2.$$

We will leave it for you to show that Eq. (34.25), together with Eq. (34.24), gives the general result

$$\mathbf{J} \cdot \mathbf{J} = j(j+1) \hbar^2. \quad (34.26)$$

Although we would think classically that the largest possible value of the z-component of \mathbf{J} is just the magnitude of \mathbf{J} —namely, $\sqrt{\mathbf{J} \cdot \mathbf{J}}$ —quantum mechanically the maximum of J_z is always a little less than that, because $j\hbar$ is always less than $\sqrt{j(j+1)} \hbar$. The angular momentum is never “completely along the z-direction.”

34-8 The magnetic energy of atoms

Now we want to talk again about the magnetic moment. We have said that in quantum mechanics the magnetic moment of a particular atomic system can be written in terms of the angular momentum by Eq. (34.6);

$$\boldsymbol{\mu} = -g \left(\frac{q_e}{2m} \right) \mathbf{J}, \quad (34.27)$$

where $-q_e$ and m are the charge and mass of the electron.

An atomic magnet placed in an external magnetic field will have an extra magnetic energy which depends on the component of its magnetic moment along the field direction. We know that

$$U_{\text{mag}} = -\boldsymbol{\mu} \cdot \mathbf{B}. \quad (34.28)$$

Choosing our z-axis along the direction of \mathbf{B} ,

$$U_{\text{mag}} = -\mu_z B. \quad (34.29)$$

Using Eq. (34.27), we have that

$$U_{\text{mag}} = g \left(\frac{q_e}{2m} \right) J_z B.$$

Quantum mechanics says that J_z can have only certain values: $j\hbar, (j-1)\hbar, \dots, -j\hbar$. Therefore, the magnetic energy of an atomic system is not arbitrary; it can have only certain values. Its maximum value, for instance, is

$$g \left(\frac{q_e}{2m} \right) j\hbar B.$$

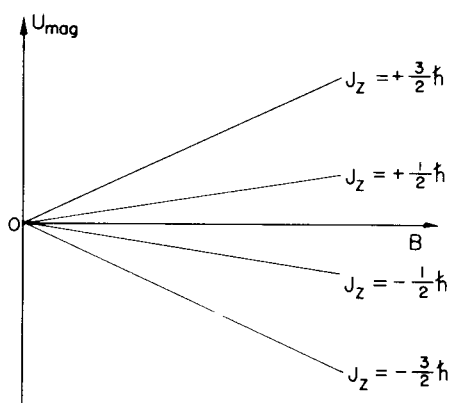


Fig. 34-5. The possible magnetic energies of an atomic system with a spin of $3/2$ in a magnetic field B .

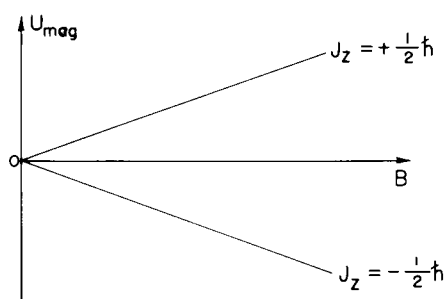


Fig. 34-6. The two possible energy states of an electron in a magnetic field B .

The quantity $q_e\hbar/2m$ is usually given the name “the Bohr magneton” and written μ_B :

$$\mu_B = \frac{q_e\hbar}{2m}.$$

The possible values of the magnetic energy are

$$U_{\text{mag}} = g\mu_B B \frac{J_z}{\hbar},$$

where J_z/\hbar takes on the possible values $j, (j-1), (j-2), \dots, (-j+1), -j$.

In other words, the energy of an atomic system is changed when it is put in a magnetic field by an amount that is proportional to the field, and proportional to J_z . We say that the energy of an atomic system is “split into $2j+1$ levels” by a magnetic field. For instance, an atom whose energy is U_0 outside a magnetic field and whose j is $3/2$, will have four possible energies when placed in a field. We can show these energies by an energy-level diagram like that drawn in Fig 34-5. Any particular atom can have only one of the four possible energies in any given field B . That is what quantum mechanics says about the behavior of an atomic system in a magnetic field.

The simplest “atomic” system is a single electron. The spin of an electron is $1/2$, so there are two possible states. $J_z = \hbar/2$ and $J_z = -\hbar/2$. For an electron at rest (no orbital motion), the spin magnetic moment has a g -value of 2, so the magnetic energy can be either $\pm\mu_B B$. The possible energies in a magnetic field are shown in Fig. 34-6. Speaking loosely we say that the electron either has its spin “up” (along the field) or “down” (opposite the field).

For systems with higher spins, there are more states. We can think that the spin is “up” or “down” or cocked at some “angle” in between, depending on the value of J_z .

We will use these quantum mechanical results to discuss the magnetic properties of materials in the next chapter.

Paramagnetism and Magnetic Resonance

35-1 Quantized magnetic states

In the last chapter we described how in quantum mechanics the angular momentum of a thing does not have an arbitrary direction, but its component along a given axis can take on only certain equally spaced, discrete values. It is a shocking and peculiar thing. You may think that perhaps we should not go into such things until your minds are more advanced and ready to accept this kind of an idea. Actually, your minds will never become more advanced—in the sense of being able to accept such a thing easily. There isn't any descriptive way of making it intelligible that isn't so subtle and advanced in its own form that it is more complicated than the thing you were trying to explain. The behavior of matter on a small scale—as we have remarked many times—is different from anything that you are used to and is very strange indeed. As we proceed with classical physics, it is a good idea to try to get a growing acquaintance with the behavior of things on a small scale, at first as a kind of experience without any deep understanding. Understanding of these matters comes very slowly, if at all. Of course, one does get better able to know what is going to happen in a quantum-mechanical situation—if that is what understanding means—but one never gets a comfortable feeling that these quantum-mechanical rules are “natural.” Of course they *are*, but they are not natural to our own experience at an ordinary level. We should explain that the attitude that we are going to take with regard to this rule about angular momentum is quite different from many of the other things we have talked about. We are not going to try to “explain” it, but we must at least *tell* you what happens; it would be dishonest to describe the magnetic properties of materials without mentioning the fact that the classical description of magnetism—of angular momentum and magnetic moments—is incorrect.

One of the most shocking and disturbing features about quantum mechanics is that if you take the angular momentum along any particular axis you find that it is always an integer or half-integer times \hbar . This is so no matter which axis you take. The subtleties involved in that curious fact—that you can take any other axis and find that the component for it is also locked to the same set of values—we will leave to a later chapter, when you will experience the delight of seeing how this apparent paradox is ultimately resolved.

We will now just accept the fact that for every atomic system there is a number j , called the *spin* of the system—which must be an integer or a half-integer—and that the component of the angular momentum along any particular axis will always have one of the following values between $+j\hbar$ and $-j\hbar$:

$$J_z = \text{one of } \left\{ \begin{array}{c} j \\ j-1 \\ j-2 \\ \vdots \\ -j+2 \\ -j+1 \\ -j \end{array} \right\} \cdot \hbar. \quad (35.1)$$

We have also mentioned that every simple atomic system has a magnetic moment which has the same direction as the angular momentum. This is true not only for atoms and nuclei but also for the fundamental particles. Each fundamental particle has its own characteristic value of j and its magnetic moment.

35-1 Quantized magnetic states

35-2 The Stern-Gerlach experiment

35-3 The Rabi molecular-beam method

35-4 The paramagnetism of bulk materials

35-5 Cooling by adiabatic demagnetization

35-6 Nuclear magnetic resonance

Review. Chapter 11, *Inside Dielectrics*

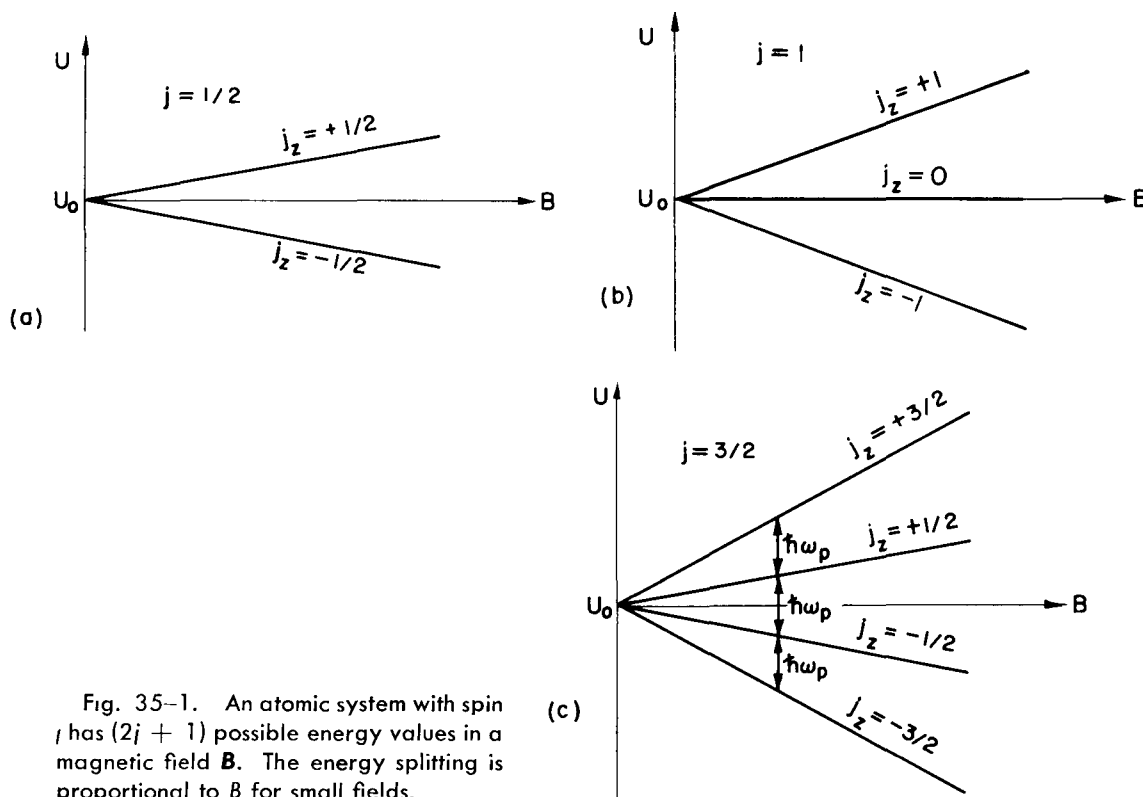


Fig. 35-1. An atomic system with spin j has $(2j + 1)$ possible energy values in a magnetic field B . The energy splitting is proportional to B for small fields.

(For some particles, both are zero.) What we mean by “the magnetic moment” in this statement is that the energy of the system in a magnetic field, say in the z -direction, can be written as $-\mu_z B$ for small magnetic fields. We must have the condition that the field should not be too great, otherwise it could disturb the internal motions of the system and the energy would not be a measure of the magnetic moment that was there before the field was turned on. But if the field is sufficiently weak, the field changes the energy by the amount

$$\Delta U = -\mu_z B, \quad (35.2)$$

with the understanding that in this equation we are to replace μ_z by

$$\mu_z = g \left(\frac{q}{2m} \right) J_z, \quad (35.3)$$

where J_z has one of the values in Eq. (35.1).

Suppose we take a system with a spin $j = 3/2$. Without a magnetic field, the system has four different possible states corresponding to the different values of J_z , all of which have exactly the same energy. But the moment we turn on the magnetic field, there is an additional energy of interaction which separates these states into four slightly different energy levels. The energies of these levels are given by a certain energy proportional to B , multiplied by \hbar times $3/2, 1/2, -1/2$, and $-3/2$ —the values of J_z . The splitting of the energy levels for atomic systems with spins of $1/2, 1$, and $3/2$ are shown in the diagrams of Fig. 35-1. (Remember that for any arrangement of electrons the magnetic moment is always directed opposite to the angular momentum.)

You will notice from the diagrams that the “center of gravity” of the energy levels is the same with and without a magnetic field. Also notice that the spacings from one level to the next are always equal for a given particle in a given magnetic field. We are going to write the energy spacing, for a given magnetic field B , as $\hbar\omega_p$ —which is just a definition of ω_p . Using Eqs. (35.2) and (35.3), we have

$$\begin{aligned} \hbar\omega_p &= g \frac{q}{2m} \hbar B \\ \text{or} \quad \omega_p &= g \frac{q}{2m} B. \end{aligned} \quad (35.4)$$

The quantity $g(q/2m)$ is just the ratio of the magnetic moment to the angular momentum—it is a property of the particle. Equation (35.4) is the same formula that we got in Chapter 34 for the angular velocity of precession in a magnetic field, for a gyroscope whose angular momentum is \mathbf{J} and whose magnetic moment is $\boldsymbol{\mu}$.

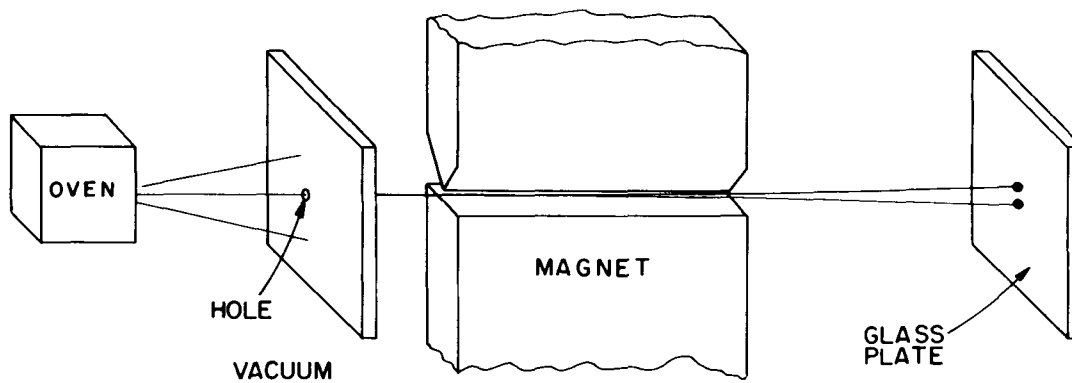


Fig. 35-2. The experiment of Stern and Gerlach.

35-2 The Stern-Gerlach experiment

The fact that the angular momentum is quantized is such a surprising thing that we will talk a little bit about it historically. It was a shock from the moment it was discovered (although it was expected theoretically). It was first observed in an experiment done in 1922 by Stern and Gerlach. If you wish, you can consider the experiment of Stern-Gerlach as a direct justification for a belief in the quantization of angular momentum. Stern and Gerlach devised an experiment for measuring the magnetic moment of individual silver atoms. They produced a beam of silver atoms by evaporating silver in a hot oven and letting some of them come out through a series of small holes. This beam was directed between the pole tips of a special magnet, as shown in Fig. 35-2. Their idea was the following. If the silver atom has a magnetic moment $\boldsymbol{\mu}$, then in a magnetic field \mathbf{B} it has an energy $-\boldsymbol{\mu}_z B$, where z is the direction of the magnetic field. In the classical theory, μ_z would be equal to the magnetic moment times the cosine of the angle between the moment and the magnetic field, so the extra energy in the field would be

$$\Delta U = -\boldsymbol{\mu} B \cos \theta. \quad (35.5)$$

Of course, as the atoms come out of the oven, their magnetic moments would point in every possible direction, so there would be all values of θ . Now if the magnetic field varies very rapidly with z —if there is a strong field gradient—then the magnetic energy will also vary with position, and there will be a force on the magnetic moments whose direction will depend on whether cosine θ is positive or negative. The atoms will be pulled up or down by a force proportional to the derivative of the magnetic energy; from the principle of virtual work,

$$F_z = -\frac{\partial U}{\partial z} = \mu \cos \theta \frac{\partial B}{\partial z}. \quad (35.6)$$

Stern and Gerlach made their magnet with a very sharp edge on one of the pole tips in order to produce a very rapid variation of the magnetic field. The beam of silver atoms was directed right along this sharp edge, so that the atoms would feel a vertical force in the inhomogeneous field. A silver atom with its magnetic moment directed horizontally would have no force on it and would go straight past the magnet. An atom whose magnetic moment was exactly vertical would have a force pulling it up toward the sharp edge of the magnet. An atom whose magnetic moment was pointed downward would feel a downward push. Thus,

as they left the magnet, the atoms would be spread out according to their vertical components of magnetic moment. In the classical theory all angles are possible, so that when the silver atoms are collected by deposition on a glass plate, one should expect a smear of silver along a vertical line. The height of the line would be proportional to the magnitude of the magnetic moment. The abject failure of classical ideas was completely revealed when Stern and Gerlach saw what actually happened. They found on the glass plate two distinct spots. The silver atoms had formed two beams.

That a beam of atoms whose spins would apparently be randomly oriented gets split up into two separate beams is most miraculous. How does the magnetic moment *know* that it is only allowed to take on certain components in the direction of the magnetic field? Well, that was really the beginning of the discovery of the quantization of angular momentum, and instead of trying to give you a theoretical explanation, we will just say that you are stuck with the result of this experiment just as the physicists of that day had to accept the result when the experiment was done. It is an *experimental fact* that the energy of an atom in a magnetic field takes on a series of individual values. For each of these values the energy is proportional to the field strength. So in a region where the field varies, the principle of virtual work tells us that the possible magnetic force on the atoms will have a set of separate values, the force is different for each state, so the beam of atoms is split into a small number of separate beams. From a measurement of the deflection of the beams, one can find the strength of the magnetic moment.

35-3 The Rabi molecular-beam method

We would now like to describe an improved apparatus for the measurement of magnetic moments which was developed by I. I. Rabi and his collaborators. In the Stern-Gerlach experiment the deflection of atoms is very small, and the measurement of the magnetic moment is not very precise. Rabi's technique permits a fantastic precision in the measurement of the magnetic moments. The method is based on the fact that the original energy of the atoms in a magnetic field is split up into a finite number of energy levels. That the energy of an atom in the magnetic field can have only certain discrete energies is really not more surprising than the fact that atoms *in general* have only certain discrete energy levels—something we mentioned often in Volume I. Why should the same thing *not* hold for atoms in a magnetic field? It does. But it is the attempt to correlate this with the idea of an *oriented magnetic moment* that brings out some of the strange implications of quantum mechanics.

When an atom has two levels which differ in energy by the amount ΔU , it can make a transition from the upper level to the lower level by emitting a light quantum of frequency ω , where

$$\hbar\omega = \Delta U \quad (35.7)$$

The same thing can happen with atoms in a magnetic field. Only then, the energy differences are so small that the frequency does not correspond to light, but to microwaves or to radiofrequencies. The transitions from the lower energy level to an upper energy level of an atom can also take place with the absorption of light or, in the case of atoms in a magnetic field, by the absorption of microwave energy. Thus if we have an atom in a magnetic field, we can cause transitions from one state to another by applying an additional electromagnetic field of the proper frequency. In other words, if we have an atom in a strong magnetic field and we “tickle” the atom with a weak varying electromagnetic field, there will be a certain probability of knocking it to another level if the frequency is near to the ω in Eq. (35.7). For an atom in a magnetic field, this frequency is just what we have earlier called ω_p , and it is given in terms of the magnetic field by Eq. (35.4). If the atom is tickled with the wrong frequency, the chance of causing a transition is very small. Thus there is a sharp *resonance* at ω_p in the probability of causing a transition. By measuring the frequency of this resonance in a known magnetic field B , we can measure the quantity $g(q/2m)$ —and hence the g -factor—with great precision.

It is interesting that one comes to the same conclusion from a classical point of view. According to the classical picture, when we place a small gyroscope with a magnetic moment μ and an angular momentum J in an external magnetic field, the gyroscope will precess about an axis parallel to the magnetic field. (See Fig 35-3.) Suppose we ask: How can we change the angle of the classical gyroscope with respect to the field—namely, with respect to the z -axis? The magnetic field produces a torque around a *horizontal* axis. Such a torque you would think is *trying* to line up the magnet with the field, but it only causes the precession. If we want to change the angle of the gyroscope with respect to the z -axis, we must exert a torque on it *about the z -axis*. If we apply a torque which goes in the same direction as the precession, the angle of the gyroscope will change to give a smaller component of J in the z -direction. In Fig. 35-3, the angle between J and the z -axis would increase. If we try to hinder the precession, J moves toward the vertical.

For our precessing atom in a uniform magnetic field, how can we apply the kind of torque we want? The answer is: with a weak magnetic field from the side. You might at first think that the direction of this magnetic field would have to rotate with the precession of the magnetic moment, so that it was always at right angles to the moment, as indicated by the field B' in Fig. 35-4(a). Such a field works very well, but an *alternating* horizontal field is almost as good. If we have a small horizontal field B' , which is always in the x -direction (plus or minus) and which oscillates with the frequency ω_p , then on each one-half cycle the torque on the magnetic moment reverses, so that it has a cumulative effect which is almost as effective as a rotating magnetic field. Classically, then, we would expect the component of the magnetic moment along the z -direction to change if we have a very weak oscillating magnetic field at a frequency which is exactly ω_p . Classically, of course, μ_z would change continuously, but in quantum mechanics the z -component of the magnetic moment cannot adjust continuously. It must jump suddenly from one value to another. We have made the comparison between the consequences of classical mechanics and quantum mechanics to give you some clue as to what might happen classically and how it is related to what actually happens in quantum mechanics. You will notice, incidentally, that the expected resonant frequency is the same in both cases.

One additional remark: From what we have said about quantum mechanics, there is no apparent reason why there couldn't also be transitions at the frequency $2\omega_p$. It happens that there isn't any analog of this in the classical case, and also it doesn't happen in the quantum theory either—at least not for the particular method of inducing the transitions that we have described. With an oscillating horizontal magnetic field, the probability that a frequency $2\omega_p$ would cause a jump of two steps at once is zero. It is only at the frequency ω_p that transitions, either upward or downward, are likely to occur.

Now we are ready to describe Rabi's method for measuring magnetic moments. We will consider here only the operation for atoms with a spin of $1/2$. A diagram of the apparatus is shown in Fig. 35-5. There is an oven which gives out a stream of neutral atoms which passes down a line of three magnets. Magnet 1

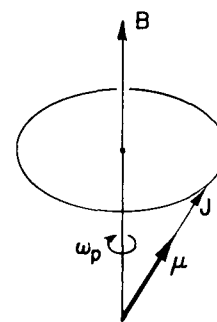


Fig. 35-3. The classical precession of an atom with the magnetic moment μ and the angular momentum J .

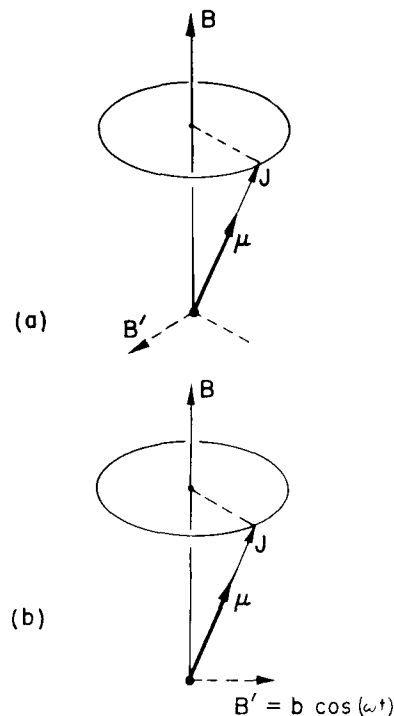


Fig. 35-4. The angle of precession of an atomic magnet can be changed by a horizontal magnetic field always at right angles to μ , as in (a), or by an oscillating field, as in (b).

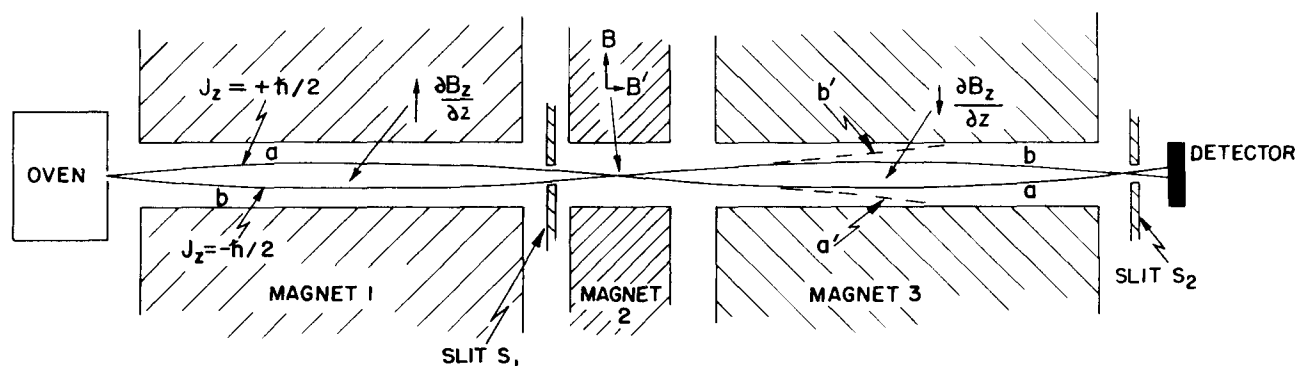


Fig. 35-5. The Rabi molecular-beam apparatus.

is just like the one in Fig. 35-2, and has a field with a strong field gradient—say, with $\partial B_z/\partial z$ positive. If the atoms have a magnetic moment, they will be deflected downward if $J_z = +\hbar/2$, or upward if $J_z = -\hbar/2$ (since for electrons μ is directed opposite to \mathbf{J}). If we consider only those atoms which can get through the slit S_1 , there are two possible trajectories, as shown. Atoms with $J_z = +\hbar/2$ must go along curve a to get through the slit, and those with $J_z = -\hbar/2$ must go along curve b . Atoms which start out from the oven along other paths will not get through the slit.

Magnet 2 has a uniform field. There are no forces on the atoms in this region, so they go straight through and enter magnet 3. Magnet 3 is just like magnet 1 but with the field *inverted*, so that $\partial B_z/\partial z$ has the opposite sign. The atoms with $J_z = +\hbar/2$ (we say “with spin up”), that felt a downward push in magnet 1, get an *upward* push in magnet 3; they continue on the path a and go through slit S_2 to a detector. The atoms with $J_z = -\hbar/2$ (“with spin down”) also have opposite forces in magnets 1 and 3 and go along the path b , which also takes them through slit S_2 to the detector.

The detector may be made in various ways, depending on the atom being measured. For example, for atoms of an alkali metal like sodium, the detector can be a thin, hot tungsten wire connected to a sensitive current meter. When sodium atoms land on the wire, they are evaporated off as Na^+ ions, leaving an electron behind. There is a current from the wire proportional to the number of sodium atoms arriving per second.

In the gap of magnet 2 there is a set of coils that produces a small horizontal magnetic field \mathbf{B}' . The coils are driven with a current which oscillates at a variable frequency ω . So between the poles of magnet 2 there is a strong, constant, vertical field \mathbf{B}_0 and a weak, oscillating, horizontal field \mathbf{B}' .

Suppose now that the frequency ω of the oscillating field is set at ω_p —the “precession” frequency of the atoms in the field \mathbf{B} . The alternating field will cause some of the atoms passing by to make transitions from one J_z to the other. An atom whose spin was initially “up” ($J_z = +\hbar/2$) may be flipped “down” ($J_z = -\hbar/2$). Now this atom has the direction of its magnetic moment reversed, so it will feel a *downward* force in magnet 3 and will move along the path a' , shown in Fig. 35-5. It will no longer get through the slit S_2 to the detector. Similarly, some of the atoms whose spins were initially down ($J_z = -\hbar/2$) will have their spins flipped up ($J_z = +\hbar/2$) as they pass through magnet 2. They will then go along the path b' and will not get to the detector.

If the oscillating field \mathbf{B}' has a frequency appreciably different from ω_p , it will not cause any spin flips, and the atoms will follow their undisturbed paths to the detector. So you can see that the “precession” frequency ω_p of the atoms in the field \mathbf{B}_0 can be found by varying the frequency ω of the field \mathbf{B}' until a decrease is observed in the current of atoms arriving at the detector. A decrease in the current will occur when ω is “in resonance” with ω_p . A plot of the detector current as a function of ω might look like the one shown in Fig. 35-6. Knowing ω_p , we can obtain the g -value of the atom.

Such atomic-beam or, as they are usually called, “molecular” beam resonance experiments are a beautiful and delicate way of measuring the magnetic properties of atomic objects. The resonance frequency ω_p can be determined with great precision—in fact, with a greater precision than we can measure the magnetic field \mathbf{B}_0 , which we must know to find g .

35-4 The paramagnetism of bulk materials

We would like now to describe the phenomenon of the paramagnetism of bulk materials. Suppose we have a substance whose atoms have permanent magnetic moments, for example a crystal like copper sulfate. In the crystal there are copper ions whose inner electron shells have a net angular momentum and a net magnetic moment. So the copper ion is an object which has a permanent magnetic moment. Let's say just a word about which atoms have magnetic moments and which ones don't. Any atom, like sodium for instance, which has an *odd* number

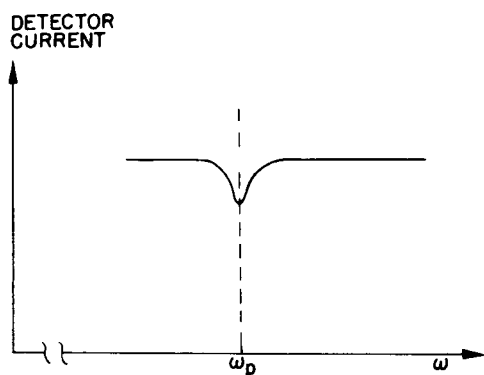


Fig. 35-6. The current of atoms in the beam decreases when $\omega = \omega_p$.

of electrons, will have a magnetic moment. Sodium has one electron in its unfilled shell. This electron gives the atom a spin and a magnetic moment. Ordinarily, however, when compounds are formed the extra electrons in the outside shell are coupled together with other electrons whose spin directions are exactly opposite, so that all the angular momenta and magnetic moments of the valence electrons usually cancel out. That's why, in general, molecules do not have a magnetic moment. Of course if you have a gas of sodium atoms, there is no such cancellation.* Also, if you have what is called in chemistry a "free radical"—an object with an odd number of valence electrons—then the bonds are not completely satisfied, and there is a net angular momentum.

In most bulk materials there is a net magnetic moment only if there are atoms present whose *inner* electron shell is not filled. Then there can be a net angular momentum and a magnetic moment. Such atoms are found in the "transition element" part of the periodic table—for instance, chromium, manganese, iron, nickel, cobalt, palladium, and platinum are elements of this kind. Also, all of the rare earth elements have unfilled inner shells and permanent magnetic moments. There are a couple of other strange things that also happen to have magnetic moments, such as liquid oxygen, but we will leave it to the chemistry department to explain the reason.

Now suppose that we have a box full of atoms or molecules with permanent moments—say a gas, or a liquid, or a crystal. We would like to know what happens if we apply an external magnetic field. With *no* magnetic field, the atoms are kicked around by the thermal motions, and the moments wind up pointing in all directions. But when there is a magnetic field, it acts to line up the little magnets; then there are more moments lying toward the field than away from it. The material is "magnetized."

We define the *magnetization* \mathbf{M} of a material as the net magnetic moment per unit volume, by which we mean the vector sum of all the atomic magnetic moments in a unit volume. If there are N atoms per unit volume and their *average* moment is $\langle \boldsymbol{\mu} \rangle_{\text{av}}$ then \mathbf{M} can be written as N times the average atomic moment:

$$\mathbf{M} = N \langle \boldsymbol{\mu} \rangle_{\text{av}} \quad (35.8)$$

The definition of \mathbf{M} corresponds to the definition of the electric polarization \mathbf{P} of Chapter 10.

The classical theory of paramagnetism is just like the theory of the dielectric constant we showed you in Chapter 11. One assumes that each of the atoms has a magnetic moment $\boldsymbol{\mu}$, which always has the same magnitude but which can point in any direction. In a field \mathbf{B} , the magnetic energy is $-\boldsymbol{\mu} \cdot \mathbf{B} = -\mu B \cos \theta$, where θ is the angle between the moment and the field. From statistical mechanics, the relative probability of having any angle is $e^{-\text{energy}/kT}$, so angles near zero are more likely than angles near π . Proceeding exactly as we did in Section 11-3, we find that for small magnetic fields \mathbf{M} is directed parallel to \mathbf{B} and has the magnitude

$$M = \frac{N\mu^2 B}{3kT}. \quad (35.9)$$

[See Eq. (11.20).] This approximate formula is correct only for $\mu B/kT$ much less than one.

We find that the induced magnetization—the magnetic moment per unit volume—is proportional to the magnetic field. This is the phenomenon of paramagnetism. You will see that the effect is stronger at lower temperatures and weaker at higher temperatures. When we put a field on a substance, it develops, for small fields, a magnetic moment proportional to the field. The ratio of \mathbf{M} to \mathbf{B} (for small fields) is called the magnetic *susceptibility*.

Now we want to look at paramagnetism from the point of view of quantum mechanics. We take first the case of an atom with a spin of $1/2$. In the absence of

* Ordinary Na vapor is mostly monatomic, although there are also some molecules of Na_2 .

a magnetic field the atoms have a certain energy, but in a magnetic field there are two possible energies, one for each value of J_z . For $J_z = +\hbar/2$, the energy is changed by the magnetic field by the amount

$$\Delta U_1 = +g \left(\frac{q_e \hbar}{2m} \right) \cdot \frac{1}{2} \cdot B. \quad (35.10)$$

(The energy shift ΔU is positive for an atom because the electron charge is negative.) For $J_z = -\hbar/2$, the energy is changed by the amount

$$\Delta U_2 = -g \left(\frac{q_e \hbar}{2m} \right) \cdot \frac{1}{2} \cdot B \quad (35.11)$$

To save writing, let's set

$$\mu_0 = g \left(\frac{q_e \hbar}{2m} \right) \cdot \frac{1}{2}; \quad (35.12)$$

then

$$\Delta U = \pm \mu_0 B. \quad (35.13)$$

The meaning of μ_0 is clear: $-\mu_0$ is the z-component of the magnetic moment in the up-spin case, and $+\mu_0$ is the z-component of the magnetic moment in the down-spin case.

Now statistical mechanics tells us that the probability that an atom is in one state or another is proportional to

$$e^{-(\text{Energy of state})/kT}.$$

With no magnetic field the two states have the same energy; so when there is equilibrium in a magnetic field, the probabilities are proportional to

$$e^{-\Delta U/kT}. \quad (35.14)$$

The number of atoms per unit volume with spin up is

$$N_{\text{up}} = a e^{-\mu_0 B/kT}, \quad (35.15)$$

and the number with spin down is

$$N_{\text{down}} = a e^{+\mu_0 B/kT}. \quad (35.16)$$

The constant a is to be determined so that

$$N_{\text{up}} + N_{\text{down}} = N, \quad (35.17)$$

the total number of atoms per unit volume. So we get that

$$a = \frac{N}{e^{+\mu_0 B/kT} + e^{-\mu_0 B/kT}}. \quad (35.18)$$

What we are interested in is the *average* magnetic moment along the z-axis. The atoms with spin up will contribute a moment of $-\mu_0$, and those with spin down will have a moment of $+\mu_0$, so the average moment is

$$\langle \mu \rangle_{\text{av}} = \frac{N_{\text{up}}(-\mu_0) + N_{\text{down}}(+\mu_0)}{N}. \quad (35.19)$$

The magnetic moment per unit volume M is then $N \langle \mu \rangle_{\text{av}}$. Using Eqs. (35.15), (35.16), and (35.17), we get that

$$M = N \mu_0 \frac{e^{+\mu_0 B/kT} - e^{-\mu_0 B/kT}}{e^{+\mu_0 B/kT} + e^{-\mu_0 B/kT}}. \quad (35.20)$$

This is the quantum-mechanical formula for M for atoms with $j = 1/2$. Incidentally, this formula can also be written somewhat more concisely in terms of the

hyperbolic tangent function:

$$M = N\mu_0 \tanh \frac{\mu_0 B}{kT}. \quad (35.21)$$

A plot of M as a function of B is given in Fig. 35.7. When B gets very large, the hyperbolic tangent approaches 1, and M approaches the limiting value $N\mu_0$. So at high fields, the magnetization *saturates*. We can see why that is; at high enough fields the moments are all lined up in the same direction. In other words, they are all in the spin-down state, and each atom contributes the moment μ_0 .

In most normal cases—say, for typical moments, room temperatures, and the fields one can normally get (like 10,000 gauss)—the ratio $\mu_0 B/kT$ is about 0.02. One must go to very low temperatures to see the saturation. For normal temperatures, we can usually replace $\tanh x$ by x , and write

$$M = \frac{N\mu_0^2 B}{kT}. \quad (35.22)$$

Just as we saw in the classical theory, M is proportional to B . In fact, the formula is almost exactly the same, except that there seems to be a factor of $1/3$ missing. But we still need to relate the μ_0 in our quantum formula to the μ that appears in the classical result, Eq. (35.9).

In the classical formula, what appears is $\mu^2 = \mu \cdot \mu$, the square of the vector magnetic moment, or

$$\mu \cdot \mu = \left(g \frac{q_e}{2m} \right)^2 J \cdot J. \quad (35.23)$$

We pointed out in the last chapter that you can very likely get the right answer from a classical calculation by replacing $J \cdot J$ by $j(j+1)\hbar^2$. In our particular example, we have $j = 1/2$, so

$$j(j+1)\hbar^2 = \frac{3}{4}\hbar^2.$$

Substituting this for $J \cdot J$ in Eq. (35.23), we get

$$\mu \cdot \mu = \left(g \frac{q_e}{2m} \right)^2 \frac{3\hbar^2}{4},$$

or in terms of μ_0 , defined in Eq. (35.12), we get

$$\mu \cdot \mu = 3\mu_0^2.$$

Substituting this for μ^2 in the classical formula, Eq. (35.9), does indeed reproduce the correct quantum formula, Eq. (35.22).

The quantum theory of paramagnetism is easily extended to atoms of any spin j . The low-field magnetization is

$$M = Ng^2 \frac{j(j+1)}{3} \frac{\mu_B^2 B}{kT}, \quad (35.24)$$

where

$$\mu_B = \frac{q_e \hbar}{2m} \quad (35.25)$$

is a combination of constants with the dimensions of a magnetic moment. Most atoms have moments of roughly this size. It is called the *Bohr magneton*. The spin magnetic moment of the electron is almost exactly one Bohr magneton.

35-5 Cooling by adiabatic demagnetization

There is a very interesting special application of paramagnetism. At very low temperatures it is possible to line up the atomic magnets in a strong field. It is then possible to get down to *extremely* low temperatures by a process called *adiabatic demagnetization*. We can take a paramagnetic salt (for example, one

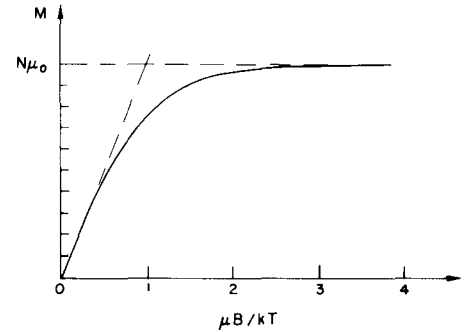


Fig. 35-7. The variation of the paramagnetic magnetization with the magnetic field strength B .

containing a number of rare-earth atoms like praseodymium-ammonium-nitrate), and start by cooling it down with liquid helium to one or two degrees absolute in a strong magnetic field. Then the factor $\mu B/kT$ is larger than 1—say more like 2 or 3. Most of the spins are lined up, and the magnetization is nearly saturated. Let's say, to make it easy, that the field is very powerful and the temperature is very low, so that nearly all the atoms are lined up. Then you isolate the salt thermally (say, by removing the liquid helium and leaving a good vacuum) and turn off the magnetic field. The temperature of the salt goes way down.

Now if you were to turn off the field *suddenly*, the jiggling and shaking of the atoms in the crystal lattice would gradually knock all the spins out of alignment. Some of them would be up and some down. But if there is no field (and disregarding the interactions between the atomic magnets, which will make only a slight error), it takes no energy to turn over the atomic magnets. They could randomize their spins without any energy change and, therefore, without any temperature change.

Suppose, however, that while the atomic magnets are being flipped over by the thermal motion there is still some magnetic field present. Then it requires some work to flip them over opposite to the field—they *must do work against the field*. This takes energy from the thermal motions and lowers the temperature. So if the strong magnetic field is not removed too rapidly, the temperature of the salt will decrease—it is cooled by the demagnetization. From the quantum-mechanical view, when the field is strong all the atoms are in the lowest state, because the odds against any being in the upper state are impossibly big. But as the field is lowered, it gets more and more likely that thermal fluctuations will knock an atom into the upper state. When that happens, the atom absorbs the energy $\Delta U = \mu_0 B$. So if the field is turned off slowly, the magnetic transitions can take energy out of the thermal vibrations of the crystal, cooling it off. It is possible in this way to go from a temperature of a few degrees absolute down to a temperature of a few thousandths of a degree.

Would you like to make something even colder than that? It turns out that Nature has provided a way. We have already mentioned that there are also magnetic moments for the atomic nuclei. Our formulas for paramagnetism work just as well for nuclei, except that the moments of nuclei are roughly a *thousand times smaller*. [They are of the order of magnitude of $qh/2m_p$, where m_p is the *proton* mass, so they are smaller by the ratio of the masses of the electron and proton.] With such magnetic moments, even at a temperature of 2°K, the factor $\mu B/kT$ is only a few parts in a thousand. But if we use the paramagnetic demagnetization process to get down to a temperature of a few thousandths of a degree, $\mu B/kT$ becomes a number near 1—at these low temperatures we can begin to saturate the nuclear moments. That is good luck, because we can then use the adiabatic demagnetization of the *nuclear* magnetism to reach still lower temperatures. Thus it is possible to do two stages of magnetic cooling. First we use adiabatic demagnetization of paramagnetic ions to reach a few thousandths of a degree. Then we use the cold paramagnetic salt to cool some material which has a strong nuclear magnetism. Finally, when we remove the magnetic field from this material, its temperature will go down to within a *millionth* of a degree of absolute zero—if we have done everything very carefully.

35-6 Nuclear magnetic resonance

We have said that atomic paramagnetism is very small and that nuclear magnetism is even a thousand times smaller. Yet it is relatively easy to observe the nuclear magnetism by the phenomenon of “nuclear magnetic resonance.” Suppose we take a substance like water, in which all of the electron spins are exactly balanced so that their net magnetic moment is zero. The molecules will still have a very, very tiny magnetic moment due to the nuclear magnetic moment of the hydrogen nuclei. Suppose we put a small sample of water in a magnetic field B . Since the protons (of the hydrogen) have a spin of 1/2, they will have two possible energy states. If the water is in thermal equilibrium, there will be slightly more

protons in the lower energy states—with their moments directed parallel to the field. There is a small net magnetic moment per unit volume. Since the proton moment is only about one-thousandth of an atomic moment, the magnetization which goes as μ^2 —using Eq. (35.22)—is only about one-millionth as strong as typical atomic paramagnetism. (That's why we have to pick a material with no atomic magnetism.) If you work it out, the difference between the number of protons with spin up and with spin down is only one part in 10^8 , so the effect is indeed very small! It can still be observed, however, in the following way.

Suppose we surround the water sample with a small coil that produces a small horizontal oscillating magnetic field. If this field oscillates at the frequency ω_p , it will induce transitions between the two energy states—just as we described for the Rabi experiment in Section 35-3. When a proton flips from an upper energy state to a lower one, it will give up the energy $\mu_z B$ which, as we have seen, is equal to $\hbar\omega_p$. If it flips from the lower energy state to the upper one, it will *absorb* the energy $\hbar\omega_p$ from the coil. Since there are slightly more protons in the lower state than in the upper one, there will be a net *absorption* of energy from the coil. Although the effect is very small, the slight energy absorption can be seen with a sensitive electronic amplifier.

Just as in the Rabi molecular-beam experiment, the energy absorption will be seen only when the oscillating field is in resonance, that is, when

$$\omega = \omega_p = g \left(\frac{q_e}{2m_p} \right) B$$

It is often more convenient to search for the resonance by varying B while keeping ω fixed. The energy absorption will evidently appear when

$$B = \frac{2m_p}{gq_e} \omega$$

A typical nuclear magnetic resonance apparatus is shown in Fig. 35-8. A high-frequency oscillator drives a small coil placed between the poles of a large electromagnet. Two small auxiliary coils around the pole tips are driven with a 60-cycle current so that the magnetic field is “wobbled” about its average value by a very small amount. As an example, say that the main current of the magnet is set to give a field of 5000 gauss, and the auxiliary coils produce a variation of ± 1 gauss about this value. If the oscillator is set at 21.2 megacycles per second, it will then be at the proton resonance each time the field sweeps through 5000 gauss [using Eq. (34.13) with $g = 5.58$ for the proton].

The circuit of the oscillator is arranged to give an additional output signal proportional to any *change* in the power being absorbed from the oscillator. This signal is fed to the vertical deflection amplifier of an oscilloscope. The horizontal sweep of the oscilloscope is triggered once during each cycle of the field-wobbling frequency. (More usually, the horizontal deflection is made to follow in proportion to the wobbling field.)

Before the water sample is placed inside the high-frequency coil, the power drawn from the oscillator is some value. (It doesn't change with the magnetic field.) When a small bottle of water is placed in the coil, however, a signal appears on the oscilloscope, as shown in the figure. We see a picture of the power being absorbed by the flipping over of the protons!

In practice, it is difficult to know how to set the main magnet to exactly 5000 gauss. What one does is to adjust the main magnet current until the resonance signal appears on the oscilloscope. It turns out that this is now the most convenient way to make an accurate measurement of the strength of a magnetic field. Of course, at some time *someone* had to measure accurately the magnetic field and frequency to determine the g -value of the proton. But now that this has been done, a proton resonance apparatus like that of the figure can be used as a “proton resonance magnetometer.”

We should say a word about the shape of the signal. If we were to wobble the magnetic field very slowly, we would expect to see a normal resonance curve. The energy absorption would read a maximum when ω_p arrived exactly at the

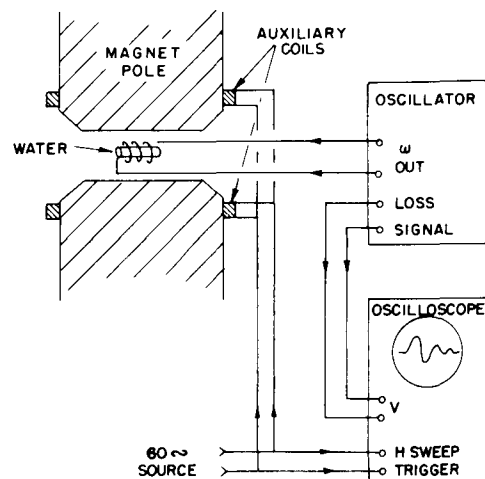


Fig 35-8. A nuclear magnetic resonance apparatus.

oscillator frequency. There would be some absorption at nearby frequencies because all the protons are not in exactly the same field—and different fields mean slightly different resonant frequencies.

One might wonder, incidentally, whether at the resonance frequency we should see any signal at all. Shouldn't we expect the high-frequency field to equalize the populations of the two states—so that there should be no signal except when the water is first put in? Not exactly, because although we are *trying* to equalize the two populations, the thermal motions on their part are trying to keep the proper ratios for the temperature T . If we sit at the resonance, the power being absorbed by the nuclei is just what is being lost to the thermal motions. There is, however, relatively little “thermal contact” between the proton magnetic moments and the atomic motions. The protons are relatively isolated down in the center of the electron distributions. So in pure water, the resonance signal is, in fact, usually too small to be seen. To increase the absorption, it is necessary to increase the “thermal contact.” This is usually done by adding a little iron oxide to the water. The iron atoms are like small magnets; as they jiggle around in their thermal dance, they make tiny jiggling magnetic fields at the protons. These varying fields “couple” the proton magnets to the atomic vibrations and tend to establish thermal equilibrium. It is through this “coupling” that protons in the higher energy states can lose their energy so that they are again capable of absorbing energy from the oscillator.

In practice the output signal of a nuclear resonance apparatus does not look like a normal resonance curve. It is usually a more complicated signal with oscillations—like the one drawn in the figure. Such signal shapes appear because of the changing fields. The explanation should be given in terms of quantum mechanics, but it can be shown that in such experiments the classical ideas of precessing moments always give the correct answer. Classically, we would say that when we arrive at resonance we start driving a lot of the precessing nuclear magnets synchronously. In so doing, we make them precess *together*. These nuclear magnets, all rotating together, will set up an induced emf in the oscillator coil at the frequency ω_p . But because the magnetic field is increasing with time, the precession frequency is increasing also, and the induced voltage is soon at a frequency a little higher than the oscillator frequency. As the induced emf goes alternately in phase and out of phase with the oscillator, the “absorbed” power goes alternately positive and negative. So on the oscilloscope we see the beat note between the proton frequency and the oscillator frequency. Because the proton frequencies are not all identical (different protons are in slightly different fields) and also possibly because of the disturbance from the iron oxide in the water, the freely precessing moments soon get out of phase, and the beat signal disappears.

These phenomena of magnetic resonance have been put to use in many ways as tools for finding out new things about matter—especially in chemistry and nuclear physics. It goes without saying that the numerical values of the magnetic moments of nuclei tell us something about their structure. In chemistry, much has been learned from the structure (or shape) of the resonances. Because of magnetic fields produced by nearby nuclei, the exact position of a nuclear resonance is shifted somewhat, depending on the environment in which any particular nucleus finds itself. Measuring these shifts helps determine which atoms are near which other ones and helps to elucidate the details of the structure of molecules. Equally important is the electron spin resonance of free radicals. Although not present to any very large extent in equilibrium, such radicals are often intermediate states of chemical reactions. A measurement of an electron spin resonance is a delicate test for the presence of free radicals and is often the key to understanding the mechanism of certain chemical reactions.

Ferromagnetism

36-1 Magnetization currents

In this chapter we will discuss some materials in which the net effect of the magnetic moments in the material is much greater than in the case of paramagnetism or diamagnetism. The phenomenon is called *ferromagnetism*. In paramagnetic and diamagnetic materials the induced magnetic moments are usually so weak that we don't have to worry about the additional fields produced by the magnetic moments. For *ferromagnetic* materials, however, the magnetic moments induced by applied magnetic fields are quite enormous and have a great effect on the fields themselves. In fact, the induced moments are so strong that they are often the dominant effect in producing the observed fields. So one of the things we will have to worry about is the mathematical theory of large induced magnetic moments. That is, of course, just a technical question. The real problem is, why are the magnetic moments so strong—how does it all work? We will come to that question in a little while.

Finding the magnetic fields of ferromagnetic materials is something like the problem of finding the electrostatic field in the presence of dielectrics. You will remember that we first described the internal properties of a dielectric in terms of a vector field \mathbf{P} , the dipole moment per unit volume. Then we figured out that the effects of this polarization are equivalent to a charge density ρ_{pol} given by the divergence of \mathbf{P} :

$$\rho_{\text{pol}} = -\nabla \cdot \mathbf{P}. \quad (36.1)$$

The total charge in any situation can be written as the sum of this polarization charge plus all other charges, whose density we write* ρ_{other} . Then the Maxwell equation which relates the divergence of \mathbf{E} to the charge density becomes

$$\nabla \cdot \mathbf{E} = \frac{\rho}{\epsilon_0} = \frac{\rho_{\text{pol}} + \rho_{\text{other}}}{\epsilon_0},$$

or

$$\nabla \cdot \mathbf{E} = -\frac{\nabla \cdot \mathbf{P}}{\epsilon_0} + \frac{\rho_{\text{other}}}{\epsilon_0}.$$

We can then pull out the polarization part of the charge and put it on the other side of the equation, to get the new law

$$\nabla \cdot (\epsilon_0 \mathbf{E} + \mathbf{P}) = \rho_{\text{other}}. \quad (36.2)$$

The new law says the divergence of the quantity $(\epsilon_0 \mathbf{E} + \mathbf{P})$ is equal to the density of the other charges.

Pulling \mathbf{E} and \mathbf{P} together as in Eq. (36.2), of course, is useful only if we know some relation between them. We have seen that the theory which relates the induced electric dipole moment to the field was a relatively complicated business and can really only be applied to certain simple situations, and even then as an approximation. We would like to remind you of one of the approximate ideas we used. To find the induced dipole moment of an atom inside a dielectric, it is necessary to know the electric field that acts on an individual atom. We made the approximation—which is not too bad in many cases—that the field on the atom

36-1 Magnetization currents

36-2 The field H

36-3 The magnetization curve

36-4 Iron-core inductances

36-5 Electromagnets

36-6 Spontaneous magnetization

Review: Chapter 10, *Dielectrics*
Chapter 17, *The Law of Induction*

* If all of the "other" charges were on conductors, ρ_{other} would be the same as our ρ_{free} of Chapter 10.

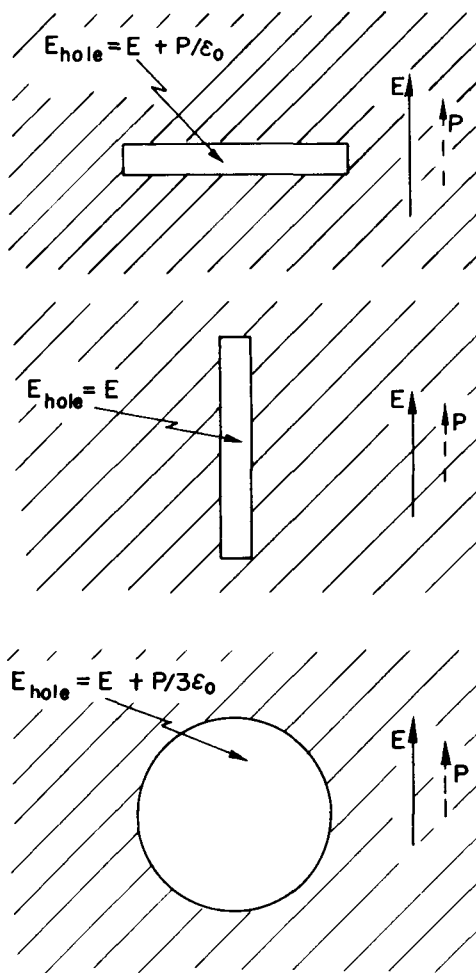


Fig. 36-1. The electric field in a cavity in a dielectric depends on the shape of the cavity.

is the same as it would be at the center of the small hole which would be left if we took out the atom (keeping the dipole moments of all the neighboring atoms the same). You will also remember that the electric field in a hole in a polarized dielectric depends on the shape of the hole. We summarize our earlier results in Fig. 36-1. For a thin, disc-shaped hole perpendicular to the polarization, the electric field in the hole is given by

$$E_{\text{hole}} = E_{\text{dielectric}} + \frac{P}{\epsilon_0},$$

which we showed by using Gauss' law. On the other hand, in a needle-shaped slot parallel to the polarization, we showed—by using the fact that the curl of \mathbf{E} is zero—that the electric fields inside and outside of the slot are the same. Finally, we found that for a spherical hole the electric field was one-third of the way between the field of the slot and the field of the disc:

$$E_{\text{hole}} = E_{\text{dielectric}} + \frac{1}{3} \frac{P}{\epsilon_0} \text{ (spherical hole)}. \quad (36.3)$$

This was the field we used in thinking about what happens to an atom inside a polarized dielectric.

Now we have to discuss the analog of all this for the case of magnetism. One simple, short-cut way of doing this is to say the \mathbf{M} , the magnetic moment per unit volume, is just like \mathbf{P} , the electric dipole moment per unit volume, and that, therefore, the negative of the divergence of \mathbf{M} is equivalent to a “magnetic charge density” ρ_m —whatever that may mean. The trouble is, of course, that there isn't any such thing as a “magnetic charge” in the physical world. As we know, the divergence of \mathbf{B} is always zero. But that does not stop us from making an artificial *analog* and writing

$$\nabla \cdot \mathbf{M} = -\rho_m, \quad (36.4)$$

where it is to be understood that ρ_m is purely mathematical. Then we could make a complete analogy with the electrostatic case and use all our old equations from electrostatics. People have often done something like that. In fact, historically, people even believed that the analogy was right. They believed that the quantity ρ_m represented the density of “magnetic poles.” These days, however, we know that the magnetization of materials comes from circulating currents within the atoms—either from the spinning electrons or from the motion of the electrons in the atom. It is therefore nicer from a physical point of view to describe things realistically in terms of the atomic currents, rather than in terms of a density of some mythical “magnetic poles.” Incidentally, these currents are sometimes called “Amperian” currents, because Ampere first suggested that the magnetism of matter came from circulating atomic currents.

The actual microscopic current density in magnetized matter is, of course, very complicated. Its value depends on where you look in the atom—it's large in some places and small in others; it goes one way in one part of the atom and the opposite way in another part (just as the microscopic electric field varies enormously inside a dielectric). In many practical problems, however, we are interested only in the fields outside of the matter or in the *average* magnetic field inside of the matter—where we mean an average taken over many, many atoms. It is only for such *macroscopic* problems that it is convenient to describe the magnetic state of the matter in terms of \mathbf{M} , the average dipole moment per unit volume. What we want to show now is that the atomic currents of magnetized matter can give rise to certain large-scale currents which are related to \mathbf{M} .

What we are going to do, then, is to separate the current density \mathbf{j} —which is the real source of the magnetic fields—into various parts: one part to describe the circulating currents of the atomic magnets, and the other parts to describe what other currents there may be. It is usually most convenient to separate the currents into three parts. In Chapter 32 we made a distinction between the currents which flow freely on conductors and the ones which are due to the back and forth motions

of the bound charges in dielectrics. In Section 32-2 we wrote

$$\mathbf{j} = \mathbf{j}_{\text{pol}} + \mathbf{j}_{\text{other}},$$

where \mathbf{j}_{pol} represented the currents from the motion of the bound charges in dielectrics and $\mathbf{j}_{\text{other}}$ took care of all other currents. Now we want to go further. We want to separate $\mathbf{j}_{\text{other}}$ into one part, \mathbf{j}_{mag} , which describes the average currents inside of magnetized materials, and an additional term which we can call \mathbf{j}_{cond} for whatever is left over. The last term will generally refer to currents in conductors, but it may also include other currents—for example the currents from charges moving freely through empty space. So we will write for the total current density:

$$\mathbf{j} = \mathbf{j}_{\text{pol}} + \mathbf{j}_{\text{mag}} + \mathbf{j}_{\text{cond}}. \quad (36.5)$$

Of course it is this total current which belongs in the Maxwell equation for the curl of \mathbf{B} :

$$c^2 \nabla \times \mathbf{B} = \frac{\mathbf{j}}{\epsilon_0} + \frac{\partial \mathbf{E}}{\partial t}. \quad (36.6)$$

Now we have to relate the current \mathbf{j}_{mag} to the magnetization vector \mathbf{M} . So that you can see where we are going, we will tell you that the result is going to be that

$$\mathbf{j}_{\text{mag}} = \nabla \times \mathbf{M}. \quad (36.7)$$

If we are given the magnetization vector \mathbf{M} everywhere in a magnetic material, the circulation current density is given by the curl of \mathbf{M} . Let's see if we can understand why this is so.

First, let's take the case of a cylindrical rod which has a uniform magnetization parallel to its axis. Physically, we know that such a uniform magnetization really means a uniform density of atomic circulating currents everywhere inside the material. Suppose we try to imagine what the actual currents would look like in a cross section of the material. We would expect to see currents something like those shown in Fig. 36-2. Each atomic current goes around and around in a little circle, with all the circulating currents going around in the same direction. Now what is the effective current of such a thing? Well, in most of the bar there is no effect at all, because right next to each current there is another current going in the opposite direction. If we imagine a small surface—but one still quite a bit larger than a single atom—such as is indicated in Fig. 36-2 by the line \overline{AB} , the net current through such a surface is zero. There is no net current anywhere inside the material. Note, however, that at the surface of the material there are atomic currents which are not cancelled by neighboring currents going the other way. At the surface there is a net current always going in the same direction around the rod. Now you see why we said earlier that a uniformly magnetized rod is equivalent to a long solenoid carrying an electric current.

How does this view fit with Eq. (36.7)? First, inside the material the magnetization \mathbf{M} is constant, so all its derivatives are zero. This agrees with our geometric picture. At the surface, however, \mathbf{M} is not really constant—it is constant up to the edge and then suddenly collapses to zero. So, right at the surface there are terrific gradients which, according to (36.7), will give a high current density. Suppose we look at what happens near the point C in Fig. 36-2. Taking the x - and y -directions as in the figure, the magnetization \mathbf{M} is in the z -direction. Writing out the components of Eq. (36.7), we have

$$\begin{aligned} \frac{\partial M_z}{\partial y} &= (j_{\text{mag}})_x, \\ -\frac{\partial M_z}{\partial x} &= (j_{\text{mag}})_y. \end{aligned} \quad (36.8)$$

At the point C , the derivative $\partial M_z / \partial y$ is zero, but $\partial M_z / \partial x$ is large and positive. Equation (36.7) says that there is a large current density in the minus y -direction. This agrees with our picture of a surface current going around the bar.

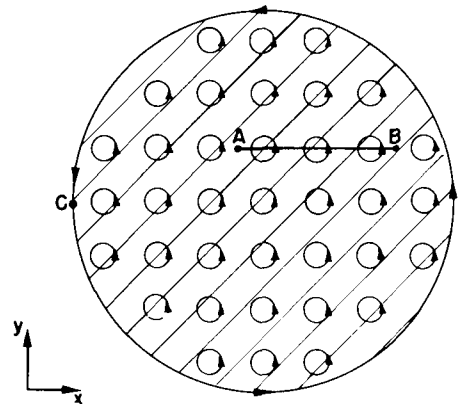


Fig. 36-2. Schematic diagram of the circulating atomic currents as seen in a cross section of an iron rod magnetized in the z -direction.

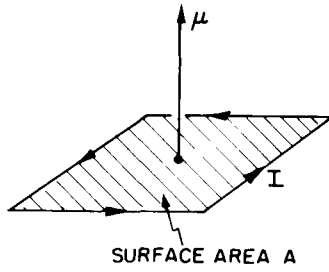


Fig. 36-3. The dipole moment μ of a current loop is IA .

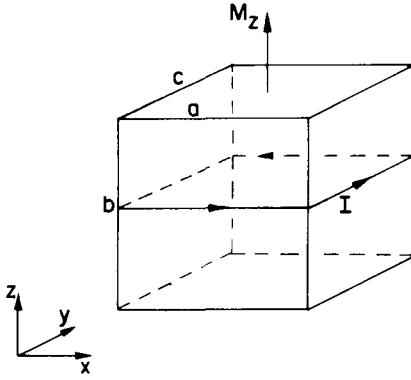
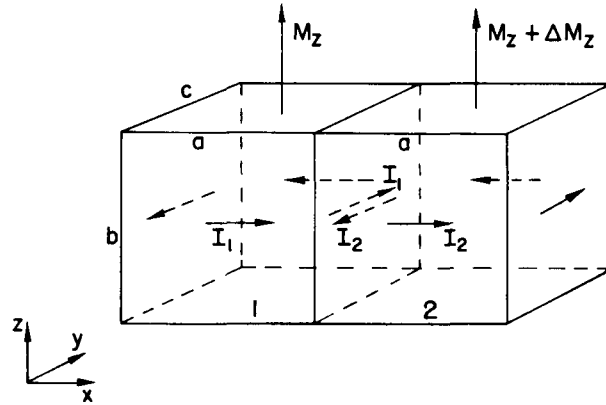


Fig. 36-4. A small magnetized block is equivalent to a circulating surface current.

Fig. 36-5. If the magnetization of two neighboring blocks is not the same, there is a net surface current in between.



Now we want to find the current density for a more complicated case in which the magnetization varies from point to point in a material. It is easy to see qualitatively that if the magnetization is different in two neighboring regions, there will not be a perfect cancellation of the circulating currents so that there will be a net current in the volume of the material. It is this effect that we want to work out quantitatively.

First, we need to recall the results of Section 14-5 that a circulating current I has a magnetic moment μ given by

$$\mu = IA, \quad (36.9)$$

where A is the area of the current loop (see Fig. 36-3). Now let's consider a small rectangular block inside of a magnetized material, as sketched in Fig. 36-4. We take the block so small that we can consider that the magnetization is uniform inside it. If this block has a magnetization M_z in the z -direction, the net effect will be the same as a surface current going around on the vertical faces, as shown. We can find the magnitude of these currents from Eq. (36.9). The total magnetic moment of the block is equal to the magnetization times the volume:

$$\mu = M_z(abc),$$

from which we get (remembering that the area of the loop is ac)

$$I = M_z b.$$

In other words, the current per unit length (vertically) on each of the vertical surfaces is equal to M_z .

Now suppose that we imagine two such little blocks next to each other, as shown in Fig. 36-5. Because block 2 is slightly displaced from block 1, it will have a slightly different vertical component of magnetization, which we call $M_z + \Delta M_z$. Now on the surface between the two blocks there will be two contributions to the total current. Block 1 will produce a current I_1 flowing in the positive y -direction, and block 2 will produce a surface current I_2 flowing in the negative y -direction. The total surface current in the positive y -direction is the sum:

$$\begin{aligned} I &= I_1 - I_2 = M_z b - (M_z + \Delta M_z)b \\ &= -\Delta M_z b. \end{aligned}$$

We can write ΔM_z as the derivative of M_z in the x -direction times the displacement from block 1 to block 2, which is just a :

$$\Delta M_z = \frac{\partial M_z}{\partial x} a.$$

The current flowing between the two blocks is then

$$I = -\frac{\partial M_z}{\partial x} ab.$$

To relate the current I to an average volume current density \mathbf{j} , we must realize that this current I is really spread over a certain cross-sectional area. If we imagine the whole volume of the material to be filled with such little blocks, one such side face (perpendicular to the x -axis) can be associated with each block.* Then we see that the area to be associated with the current I is just the area ab of one of the front faces. We get the result

$$j_y = \frac{I}{ab} = -\frac{\partial M_z}{\partial x}.$$

We have at least the beginning of the curl of \mathbf{M} .

There should be another term in j_y from the variation of the x -component of the magnetization with z . This contribution to \mathbf{j} will come from the surface between two little blocks stacked one on top of the other, as shown in Fig. 36-6. Using the same arguments we have just made, you can show that this surface will contribute to j_y the amount $\partial M_x / \partial z$. These are the only surfaces which can contribute to the y -component of the current so we have that the total current density in the y -direction is

$$j_y = \frac{\partial M_x}{\partial z} - \frac{\partial M_z}{\partial x}.$$

Working out the currents on the remaining faces of a cube—or using the fact that our z -direction is completely arbitrary—we can conclude that the vector current density is indeed given by the equation

$$\mathbf{j} = \nabla \times \mathbf{M}.$$

So if we choose to describe the magnetic situation in matter in terms of the average magnetic moment per unit volume \mathbf{M} , we find that the circulating atomic currents are equivalent to an average current density in matter given by Eq. (36.7). If the material is also a dielectric, there may be, in addition, a polarization current $\mathbf{j}_{\text{pol}} = \partial \mathbf{P} / \partial t$. And if the material is also a conductor, we may have a conduction current \mathbf{j}_{cond} as well. We can write the total current as

$$\mathbf{j} = \mathbf{j}_{\text{cond}} + \nabla \times \mathbf{M} + \frac{\partial \mathbf{P}}{\partial t}. \quad (36.10)$$

36-2 The field \mathbf{H}

Next, we want to insert the current as written in Eq. (36.10) into Maxwell's equations. We get

$$c^2 \nabla \times \mathbf{B} = \frac{\mathbf{j}}{\epsilon_0} + \frac{\partial \mathbf{E}}{\partial t} = \frac{1}{\epsilon_0} \left(\mathbf{j}_{\text{cond}} + \nabla \times \mathbf{M} + \frac{\partial \mathbf{P}}{\partial t} \right) + \frac{\partial \mathbf{E}}{\partial t}.$$

We can move the term in \mathbf{M} to the left-hand side:

$$c^2 \nabla \times \left(\mathbf{B} - \frac{\mathbf{M}}{\epsilon_0 c^2} \right) = \frac{\mathbf{j}_{\text{cond}}}{\epsilon_0} + \frac{\partial}{\partial t} \left(\mathbf{E} + \frac{\mathbf{P}}{\epsilon_0} \right). \quad (36.11)$$

As we remarked in Chapter 32, many people like to write $(\mathbf{E} + \mathbf{P}/\epsilon_0)$ as a new vector field \mathbf{D}/ϵ_0 . Similarly, it is often convenient to write $(\mathbf{B} - \mathbf{M}/\epsilon_0 c^2)$ as a single vector field. We choose to define a new vector field \mathbf{H} by

$$\mathbf{H} = \mathbf{B} - \frac{\mathbf{M}}{\epsilon_0 c^2}. \quad (36.12)$$

Then Eq. (36.11) becomes

$$\epsilon_0 c^2 \nabla \times \mathbf{H} = \mathbf{j}_{\text{cond}} + \frac{\partial \mathbf{D}}{\partial t}. \quad (36.13)$$

It looks simple, but all the complexity is just hidden in the letters \mathbf{D} and \mathbf{H} .

* Or, if you prefer, the current I in each face should be split 50-50 with the blocks on the two sides.

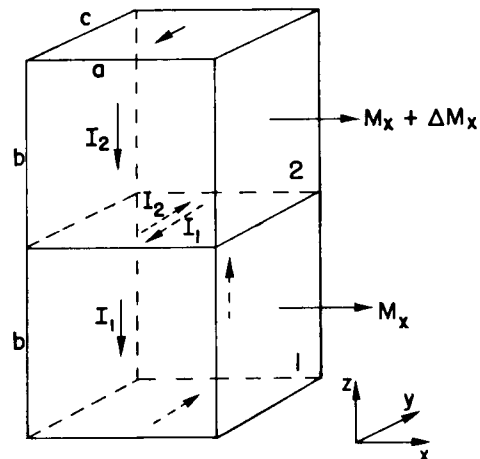


Fig. 36-6. Two blocks, one above the other, may also contribute to j_y .

Now we have to give you a warning. Most people who use the mks units have chosen to use a different definition of H . Calling *their* field H' (of course, they still call it H without the prime), it is defined by

$$H' = \epsilon_0 c^2 B - M. \quad (36.14)$$

(Also, they usually write $\epsilon_0 c^2$ as a new number $1/\mu_0$; then they have one more constant to keep track of!) With this definition, Eq. (36.13) looks even simpler:

$$\nabla \times H' = j_{\text{cond}} + \frac{\partial D}{\partial t}. \quad (36.15)$$

But the difficulties with this definition of H' are, first, that it doesn't agree with the definition of people who don't use the mks units, and second, that it makes H' and B have different units. We think it is more convenient for H to have the same units as B —rather than the units of M , as H' does. But if you are going to be an engineer and work on the design of transformers, magnets, and such, you will have to watch out. You will find many books which use for H the definition of Eq. (36.14) rather than our definition of Eq. (36.12), and many other books—especially handbooks about magnetic materials—that relate B and H the way we have done. You'll have to be careful to figure out which convention they are using.

One way to tell is by the units they use. Remember that in the mks system, B —and therefore *our* H —are measured with the unit: one weber per square meter, equal to 10,000 gauss. In the mks system, a magnetic moment (a current times an area) has the unit: one ampere-meter². The magnetization M , then, has the unit: one ampere *per* meter. For H' the units are the same as for M . You can see that this also agrees with Eq. (36.15), since ∇ has the dimensions of one over a length. People who are working with electromagnets also get in the habit of calling the unit of H (with the H' definition) “one ampere *turn* per meter”—thinking of the turns of wire on a winding. But a “turn” is really a dimensionless number, so that doesn't need to confuse you. Since our H is equal to $H'/\epsilon_0 c^2$, if you are using the mks system, H (in webers/meter²) is equal to $4\pi \times 10^{-7}$ times H' (in amperes per meter). It is perhaps more convenient to remember that H (in gauss) = 0.0126 H' (in amp/meter).

There is one more horrible thing. Many people who use *our* definition of H have decided to call the units of H and B by *different names*! Even though they have the same dimensions, they call the unit of B one *gauss*, and the unit of H one *oersted* (after Gauss and Oersted, of course). So, in many books you will find graphs with B plotted in gauss and H in oersteds. They are really the same unit— 10^{-4} of the mks unit. We have summarized the confusion about magnetic units in Table 36-1.

Table 36-1

Units of magnetic quantities

$[B]$	= weber/meter ² = 10^4 gauss
$[H]$	= weber/meter ² = 10^4 gauss or 10^4 oersted
$[M]$	= ampere/meter
$[H']$	= ampere/meter

Convenient conversions

B (gauss)	= $10^4 B$ (weber/meter ²)
H (gauss)	= H (oersted)
	= 0.0126 H' (amp/meter)

36-3 The magnetization curve

Now we will look at some simple situations in which the magnetic field is constant, or in which the fields change slowly enough that we can neglect $\partial D/\partial t$ in comparison with j_{cond} . Then the fields obey the equations

$$\nabla \cdot B = 0, \quad (36.16)$$

$$\nabla \times H = j_{\text{cond}}/\epsilon_0 c^2, \quad (36.17)$$

$$H = B - M/\epsilon_0 c^2. \quad (36.18)$$

Suppose we have a torus (a donut) of iron wrapped with a coil of copper wire, as shown in Fig. 36-7(a). A current I flows in the wire. What is the magnetic field? The magnetic field will be mainly inside the iron; there, the lines of B will be circles, as drawn in Fig. 36-7(b). Since the flux of B is continuous, its divergence is zero, and Eq. (36.16) is satisfied. Next, we write Eq. (36.17) in another form by

integrating around the closed loop Γ drawn in Fig. 36-7(b). From Stokes's theorem, we have that

$$\oint_{\Gamma} \mathbf{H} \cdot d\mathbf{s} = \frac{1}{\epsilon_0 c^2} \int_S \mathbf{j}_{\text{cond}} \cdot \mathbf{n} da, \quad (36.19)$$

where the integral of \mathbf{j} is to be carried out over any surface S bounded by Γ . This surface is cut once by each turn of the winding. Each turn contributes the current I to the integral, and, if there are N turns in all, the integral is NI . From the symmetry of our problem, \mathbf{B} is the same all around the curve Γ ; if we assume that the magnetization, and therefore, the field \mathbf{H} is also constant along Γ , Eq. (36.19) becomes

$$Hl = \frac{NI}{\epsilon_0 c^2},$$

where l is the length of the curve Γ . So,

$$H = \frac{1}{\epsilon_0 c^2} \frac{NI}{l}. \quad (36.20)$$

It is because \mathbf{H} is directly proportional to the magnetizing current in cases like this one that \mathbf{H} is sometimes called the *magnetizing field*.

Now all we need is an equation which relates \mathbf{H} to \mathbf{B} . But there isn't any such equation! There is, of course, Eq. (36.18), but it is no help because there is no direct relation between \mathbf{M} and \mathbf{B} for a ferromagnetic material like iron. The magnetization \mathbf{M} depends on the whole past history of the iron, and not only on what \mathbf{B} is at the moment.

All is not lost, though. We can get solutions in certain simple cases. If we start out with unmagnetized iron—let's say with iron that has been annealed at high temperatures—then in the simple geometry of the torus, all the iron will have the same magnetic history. Then we can say something about \mathbf{M} —and therefore about the relation between \mathbf{B} and \mathbf{H} —from experimental measurements. The field \mathbf{B} in the torus is, from Eq. (36.20), given as a constant times the current I in the winding. The field \mathbf{B} can be measured by integrating over time the emf in the coil (or in an extra coil wound over the magnetizing coil shown in the figure). This emf is equal to the rate of change of the flux of \mathbf{B} , so the integral of the emf with time is equal to \mathbf{B} times the cross-sectional area of the torus.

Figure 36-8 shows the relation between \mathbf{B} and \mathbf{H} , observed with a torus of soft iron. When the current is first turned on, \mathbf{B} increases with increasing \mathbf{H} along the curve a . Note the different scales on \mathbf{B} and \mathbf{H} ; initially, it takes only a relatively small \mathbf{H} to make a large \mathbf{B} . Why is \mathbf{B} so much larger with the iron than it would be with air? Because there is a large magnetization \mathbf{M} which is equivalent to a large surface current on the iron—the field \mathbf{B} comes from the *sum* of this current and the conduction current in the winding. Why \mathbf{M} should be so large, we will discuss later.

At higher values of \mathbf{H} , the magnetization curve levels off. We say that the iron *saturates*. With the scales of our figure, the curve appears to become horizontal. Actually, it continues to rise slightly—for large fields, \mathbf{B} becomes proportional to \mathbf{H} , and with a unit slope. There is no further increase of \mathbf{M} . Incidentally, we should point out that if the torus were made of some nonmagnetic material, \mathbf{M} would be zero and \mathbf{B} would equal \mathbf{H} for all fields.

The first thing we notice is that curve a in Fig. 36-8—which is the so-called *magnetization curve*—is highly nonlinear. But it's worse than that. If, after reaching saturation, we decrease the current in the coil to bring \mathbf{H} back to zero, the magnetic field \mathbf{B} falls along curve b . When \mathbf{H} reaches zero, there is still some \mathbf{B} left. Even with no magnetizing current there is a magnetic field in the iron—it has become permanently magnetized. If we now turn on a *negative* current in the coil, the \mathbf{B} - \mathbf{H} curve continues along b until the iron is saturated in the negative direction. If we then bring the current back to zero again, \mathbf{B} goes along curve c . If we alternate the current between large positive and negative values, the \mathbf{B} - \mathbf{H} curve goes back and forth along very nearly the curves b and c . If we vary \mathbf{H} in some arbitrary

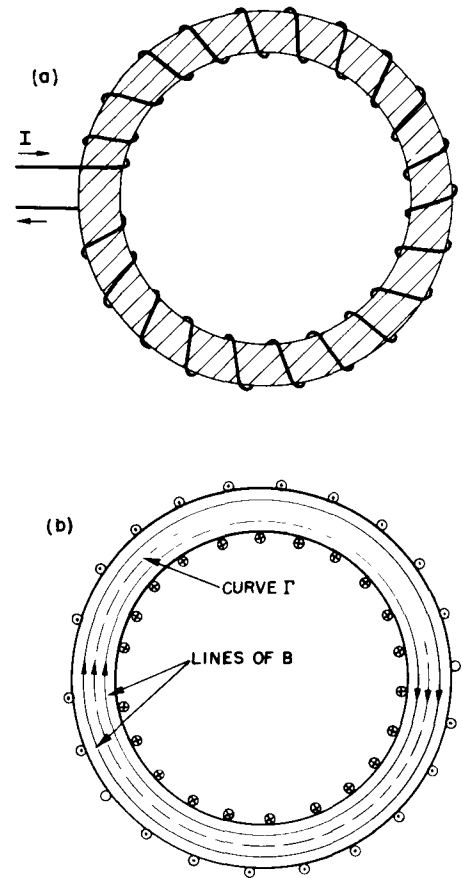


Fig. 36-7. (a) A torus of iron wound with a coil of insulated wire. (b) Cross section of torus showing field lines.

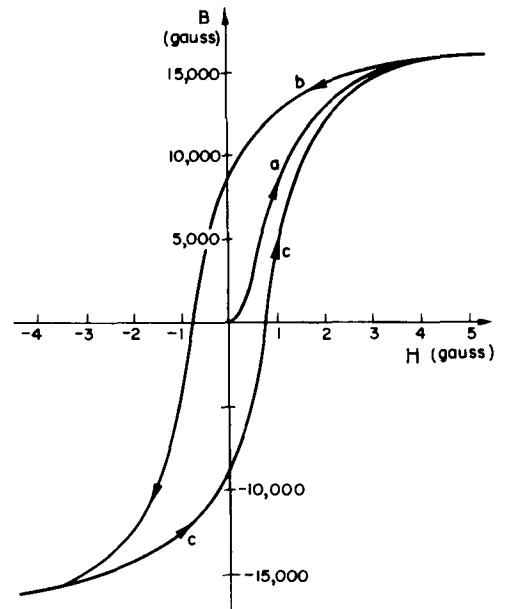


Fig. 36-8. Typical magnetization and hysteresis curves for soft iron.

way, however, we can get more complicated curves which will, in general, lie somewhere between the curves b and c . The loop made by repeated oscillation of the fields is called a *hysteresis* loop of the iron.

We see then that we cannot write a functional relationship like $B = f(H)$, because the value of B at any instant depends not only on what H is at that time, but on its whole past history. Naturally, the magnetization and hysteresis curves are different for different substances. The shape of the curves depends critically on the chemical composition of the material, and also on the details of its preparation and subsequent physical treatment. We will discuss some of the physical explanations for these complications in the next chapter.

36-4 Iron-core inductances

One of the most important applications of magnetic materials is in electrical circuits—for example, in transformers, electric motors, and so on. One reason is that with iron we can control where the magnetic fields go, and also get much larger fields for a given electric current. For example, the typical “toroidal” inductance is made very much like the object shown in Fig. 36-7. For a given inductance, it can be much smaller in volume and use much less copper than an equivalent “air-core” inductance. For a given inductance, we get a much smaller resistance in the winding, so the inductance is more nearly “ideal”—particularly for low frequencies. It is very easy to understand, qualitatively, how such an inductance works. If I is the current in the winding, then the field H which is produced in the inside is proportional to I —as given by Eq. (36.20). The voltage \mathcal{V} across the terminals is related to the magnetic field B . Neglecting the resistance of the winding, the voltage \mathcal{V} is proportional to $\partial B/\partial t$. The inductance \mathcal{L} , which is the ratio of \mathcal{V} to dI/dt (see Section 17-7), thus involves the relation between B and H in the iron. Since the B is so much bigger than the H , we get a large factor in the inductance. Physically, what happens is that a small current in the coil, which would ordinarily produce a small magnetic field, causes the little “slave” magnets in the iron to line up and produce a tremendously greater “magnetic” current than the external current in the winding. It is as if we had a lot more current going through the coil than we really have. When we reverse the current, all the little magnets flip over—all those internal currents reverse—and we get a much higher induced emf than we would get without the iron. If we want to calculate the inductance, we can do so through the energy—as described in Section 17-8. The *rate* at which energy is delivered from the current source is $I\mathcal{V}$. The voltage \mathcal{V} is the cross-sectional area A of the core, times N , times dB/dt . From Eq. (36.20), $I = (\epsilon_0 c^2 l/N)H$. So we have

$$\frac{dU}{dt} = \mathcal{V}I = (\epsilon_0 c^2 l A) H \frac{dB}{dt}.$$

Integrating over time, we have

$$U = (\epsilon_0 c^2 l A) \int H dB. \quad (36.21)$$

Notice that lA is the volume of the torus, so we have shown that the energy density $u = U/\text{vol}$ in a magnetic material is given by

$$u = \epsilon_0 c^2 \int H dB. \quad (36.22)$$

An interesting feature is involved here. When we use alternating currents, the iron is driven around a hysteresis loop. Since B is not a single-valued function of H , the integral of $\int H dB$ around one complete cycle is *not* equal to zero. It is the area enclosed inside the hysteresis curve. Thus, the driving source delivers a certain net energy each cycle—an energy proportional to the area inside the hysteresis loop. And that energy is “lost.” It is lost from the electromagnetic goings on, but turns up as heat in the iron. It is called the *hysteresis loss*. To keep such energy losses small, we would like the hysteresis loop to be as narrow as

possible. One way to decrease the area of the loop is to reduce the maximum field that is reached during each cycle. For smaller maximum fields, we get a hysteresis curve like the one shown in Fig. 36-9. Also, special materials are designed to have a very narrow loop. The so-called *transformer irons*—which are iron alloys with a small amount of silicon—have been developed to have this property.

When an inductance is run over a small hysteresis loop, the relationship between B and H can be approximated by a linear equation. People usually write

$$B = \mu H. \quad (36.23)$$

The constant μ is *not* the magnetic moment we have used before. It is called the *permeability* of the iron. (It is also sometimes called the “*relative permeability*.”) The permeability of ordinary irons is typically several thousand. There are special alloys alike “supermalloy” which can have permeabilities as high as a million.

If we use the approximation that $B = \mu H$ in Eq. (36.21), we can write the energy in a toroidal inductance as

$$U = (\epsilon_0 c^2 l A) \mu \int H dH = (\epsilon_0 c^2 l A) \frac{\mu H^2}{2}. \quad (36.24)$$

So the energy density is approximately

$$u \approx \frac{\epsilon_0 c^2}{2} \mu H^2.$$

We can now set the energy of Eq. (36.24) equal to the energy $\mathcal{L} I^2 / 2$ of an inductance, and solve for \mathcal{L} . We get

$$\mathcal{L} = (\epsilon_0 c^2 l A) \mu \left(\frac{H}{I} \right)^2.$$

Using H/I from Eq. (36.20), we have

$$\mathcal{L} = \frac{\mu N^2 A}{\epsilon_0 c^2 l}. \quad (36.25)$$

The inductance is proportional to μ . If you want inductances for such things as audio amplifiers, you will try to operate them on a hysteresis loop where the B - H relationship is as linear as possible. (You will remember that we spoke in Chapter 50, Vol. I, about the generation of harmonics in nonlinear systems.) For such purposes, Eq. (36.23) is a useful approximation. On the other hand, if you *want* to generate harmonics, you may use an inductance which is intentionally operated in a highly nonlinear way. Then you will have to use the complete B - H curves, and analyze what happens by graphical or numerical methods.

A “transformer” is often made by putting two coils on the same torus—or *core*—of a magnetic material. (For the larger transformers, the core is made with rectangular proportions for convenience.) Then a varying current in the “primary” winding causes the magnetic field in the core to change, which induces an emf in the “secondary” winding. Since the flux through *each turn* of both windings is the same, the emf’s in the two windings are in the same ratio as the number of turns on each. A voltage applied to the primary is transformed to a different voltage at the secondary. Since a certain *net* current around the core is needed to produce the required change in the magnetic field, the *algebraic* sum of the currents in the two windings will be fixed and equal to the required “magnetizing” current. If the current drawn from the secondary increases, the primary current must increase in proportion—there is a “transformation” of currents as well as voltage.

36-5 Electromagnets

Now let’s discuss a practical situation which is a little more complicated. Suppose we have an electromagnet of the rather standard form shown in Fig. 36-10—there is a “C-shaped” yoke of iron, with a coil of many turns of wire wrapped around the yoke. What is the magnetic field B in the gap?

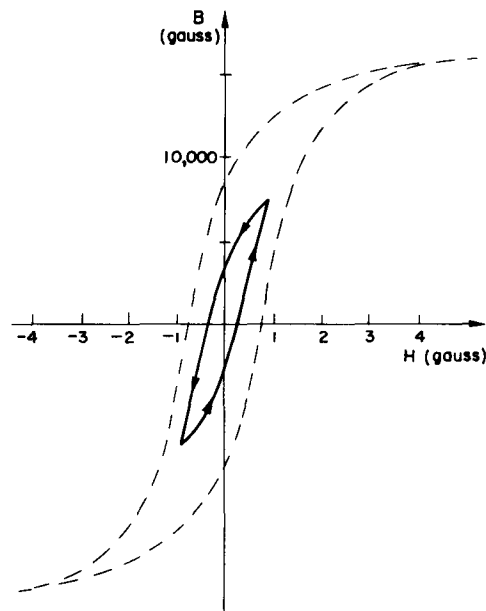


Fig. 36-9. A hysteresis loop that doesn't reach saturation.

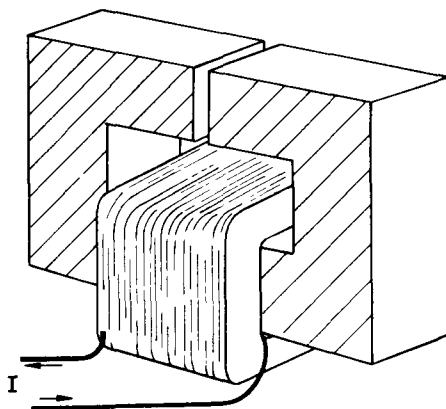


Fig. 36-10. An electromagnet.

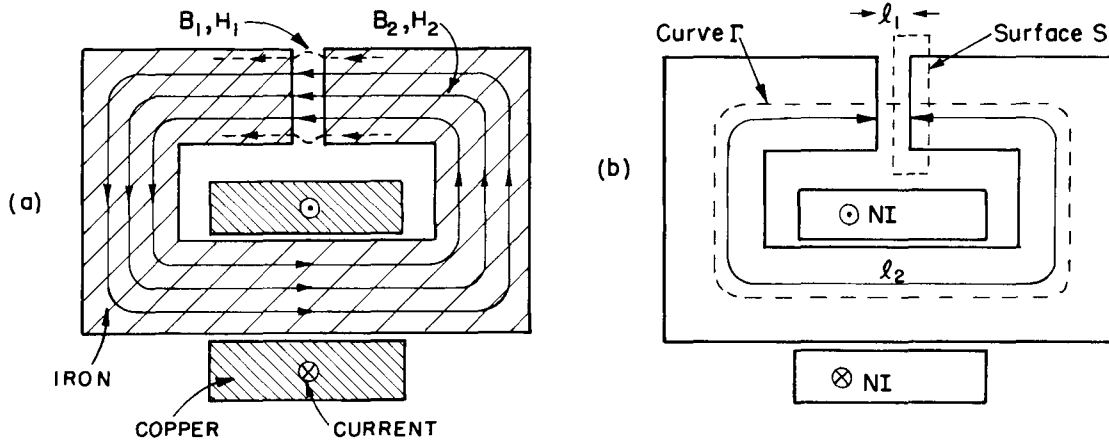


Fig. 36-11. Cross section of an electromagnet.

If the gap thickness is small compared with all the other dimensions, we can, as a first approximation, assume that the lines of \mathbf{B} will go around through the loop, just as they did in the torus. They will look more or less as shown in Fig. 36-11(a). They tend to spread out somewhat in the gap, but if the gap is narrow, this will be a small effect. It is a fair approximation to assume that the flux of \mathbf{B} through any cross section of the yoke is a constant. If the yoke has a uniform cross-sectional area—and if we neglect any edge effects at the gaps or at the corners—we can say that \mathbf{B} is uniform around the yoke.

Also, \mathbf{B} will have the same value in the gap. This follows from Eq. (36.16). Imagine the closed surface S , shown in Fig. 36-11(b), which has one face in the gap and the other in the iron. The total flux of \mathbf{B} out of this surface must be zero. Calling B_1 the field in the gap and B_2 the field in the iron, we have that

$$B_1 A_1 - B_2 A_2 = 0.$$

Since $A_1 = A_2$ (to our approximation), it follows that $B_1 = B_2$.

Now let's look at H . We can again use Eq. (36.19), taking the line integral around the curve Γ in Fig. 36-11(b). As before, the right-hand side is NI , the number of turns times the current. Now, however, H will be different in the iron and in the air. Calling H_2 the field in the iron and l_2 the path length around the yoke, this part of the curve will contribute the amount $H_2 l_2$ to the integral. Calling H_1 the field in the gap and l_1 the gap thickness, we get the contribution $H_1 l_1$ from the gap. We have that

$$H_1 l_1 + H_2 l_2 = \frac{NI}{\epsilon_0 c^2}. \quad (36.26)$$

Now we know something else: that in the air gap, the magnetization is negligible, so that $B_1 = H_1$. Since $B_1 = B_2$, Eq. (36.26) becomes

$$B_2 l_1 + H_2 l_2 = \frac{NI}{\epsilon_0 c^2}. \quad (36.27)$$

We still have two unknowns. To find B_2 and H_2 , we need another relationship—namely, the one which relates B to H in the iron.

If we can make the approximation that $B_2 = \mu H_2$, we can solve the equation algebraically. However, let's do the general case, in which the magnetization curve of the iron is one like that shown in Fig. 36-8. What we want is the simultaneous solution of this functional relationship together with Eq. (36.27). We can find it by plotting a graph of Eq. (36.27) on the same graph with the magnetization curve, as is done in Fig. 36-12. Where the two curves intersect, we have our solution.

For a given current I , the function (36.27) is the straight line marked $I > 0$ in Fig. 36-12. The line intersects the H -axis ($B_2 = 0$) at $H_2 = NI/\epsilon_0 c^2 l_2$, and the slope is $-l_2/l_1$. Different currents just shift the line horizontally. From Fig.

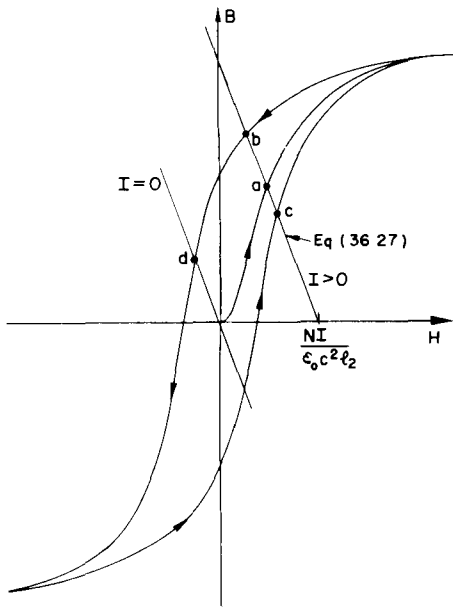


Fig. 36-12. Solving for the field in an electromagnet.

36-12, we see that for a given current there are several different solutions, depending on how you got there. If you have just built the magnet and turned the current up to I , the field B_2 (which is also B_1) will have the value given by point a . If you have run the current to some very high value and come down to I , the field will be given by point b . Or, if you have just had a high negative current in the magnet and then come up to I , the field is the one at point c . The field in the gap will depend on what you have done in the past.

When the current in the magnet is zero, the relation between B_2 and H_2 in Eq. (36.27) is shown by the line marked $I = 0$ in the figure. There are still various possible solutions. If you have first saturated the iron, there may be a considerable residual field in the magnet as given by point d . You can take the coil off, and you have a permanent magnet. You can see that for a good permanent magnet, you would want a material with a *wide* hysteresis loop. Special alloys, such as Alnico V, have very wide loops.

36-6 Spontaneous magnetization

We now turn to the question of why it is that in ferromagnetic materials a small magnetic field produces such a large magnetization. The magnetization of ferromagnetic materials like iron and nickel comes from the magnetic moment of the electrons in the inner shell of the atom. Each electron has a magnetic moment μ equal to $q/2m$ times its g -factor, times its angular momentum \mathbf{J} . For a single electron with no net orbital motion, $g = 2$, and the component of \mathbf{J} in any direction—say the z -direction—is $\pm \hbar/2$, so the component of μ along the z -axis is

$$\mu_z = \frac{q\hbar}{2m} = 0.928 \times 10^{-23} \text{ amp}\cdot\text{m}^2. \quad (36.28)$$

In an iron atom, there are actually two electrons that contribute to the ferromagnetism, so to keep the discussion simpler we will talk about nickel, which is ferromagnetic like iron but which has only one electron in the inner shell. (It is easy to extend the arguments to iron.)

Now the point is that in the presence of an external field \mathbf{B} , the atomic magnets tend to line up with the field, but are knocked about by thermal motions just as we described for paramagnetic materials. In the last chapter we found out that the balance between a magnetic field trying to line up the atomic magnets and the thermal motions trying to derange them produced the result that the mean magnetic moment per unit volume will end up as

$$M = N\mu \tanh \frac{\mu B_a}{kT}. \quad (36.29)$$

By B_a we mean the field acting at the atom, and kT is the Boltzmann energy. In the theory of paramagnetism we used for B_a just B itself, neglecting the part of the field at any given atom contributed by the atoms nearby. In the ferromagnetic case, there is a complication. We shouldn't use the average field in the iron for the B_a acting on an individual atom. Instead, we must do as we did in the case of dielectrics—we have to find the *local* field acting at a single atom. For an exact calculation we should add up the fields at the atom in question contributed by all of the other atoms in the crystal lattice. But as we did for dielectrics, we will make the approximation that the field at an atom is the same as we would find in a small spherical hole in the material—assuming that the moments of the atoms in the neighborhood are not changed by the presence of the hole.

Following the arguments we made in Chapter 11, we might think that we could write

$$B_{\text{hole}} = B + \frac{1}{3} \frac{M}{\epsilon_0 c^2} \quad (\text{wrong!}).$$

But that is not right. We *can*, however, make use of the results of Chapter 11 if we make a careful comparison of the equations of Chapter 11 with the equations

for ferromagnetism in this chapter. Let's put together the corresponding equations. For regions where there are no conduction currents or charges we have:

Electrostatics	Static ferromagnetism
$\nabla \cdot \left(\mathbf{E} + \frac{\mathbf{P}}{\epsilon_0} \right) = 0$	$\nabla \cdot \mathbf{B} = 0$
$\nabla \times \mathbf{E} = 0$	$\nabla \times \left(\mathbf{B} - \frac{\mathbf{M}}{\epsilon_0 c^2} \right) = 0$

(36.30)

These two sets of equations can be thought of as analogous if we make the following *purely mathematical* correspondences:

$$\mathbf{E} \rightarrow \mathbf{B} - \frac{\mathbf{M}}{\epsilon_0 c^2}, \quad \mathbf{E} + \frac{\mathbf{P}}{\epsilon_0} \rightarrow \mathbf{B}.$$

This is the same as making the analogy

$$\mathbf{E} \rightarrow \mathbf{H}, \quad \mathbf{P} \rightarrow \mathbf{M}/c^2. \quad (36.31)$$

In other words, if we write the equations of ferromagnetism as

$$\begin{aligned} \nabla \cdot \left(\mathbf{H} + \frac{\mathbf{M}}{\epsilon_0 c^2} \right) &= 0, \\ \nabla \times \mathbf{H} &= 0, \end{aligned} \quad (36.32)$$

they *look like* the equations of electrostatics.

This purely algebraic correspondence has led to some confusion in the past. People tended to think that \mathbf{H} was “*the* magnetic field.” But, as we have seen, \mathbf{B} and \mathbf{E} are physically the fundamental fields, and \mathbf{H} is a derived idea. So although the *equations* are analogous, the *physics* is not analogous. However, that doesn't need to stop us from using the principle that the same equations have the same solutions.

We can use our earlier results for the electric field inside of holes of various shapes in dielectrics—summarized in Fig. 36-1—to find the field \mathbf{H} inside of corresponding holes. Knowing \mathbf{H} , we can determine \mathbf{B} . For instance (using the results we summarized in Section 1), the field \mathbf{H} in a needle-shaped hole parallel to \mathbf{M} is the same as the \mathbf{H} in the material,

$$\mathbf{H}_{\text{hole}} = \mathbf{H}_{\text{material}}.$$

But since \mathbf{M} in the hole is zero, we have

$$\mathbf{B}_{\text{hole}} = \mathbf{B}_{\text{material}} - \frac{\mathbf{M}}{\epsilon_0 c^2}. \quad (36.33)$$

On the other hand, for a disc-shaped hole, perpendicular to \mathbf{M} , we have

$$\mathbf{E}_{\text{hole}} = \mathbf{E}_{\text{dielectric}} + \frac{\mathbf{P}}{\epsilon_0},$$

which translates into

$$\mathbf{H}_{\text{hole}} = \mathbf{H}_{\text{material}} + \frac{\mathbf{M}}{\epsilon_0 c^2}.$$

Or, in terms of \mathbf{B} ,

$$\mathbf{B}_{\text{hole}} = \mathbf{B}_{\text{material}}. \quad (36.34)$$

Finally, for a spherical hole, by making our analogy with Eq. (36.3) we would have

$$\mathbf{H}_{\text{hole}} = \mathbf{H}_{\text{material}} + \frac{\mathbf{M}}{3\epsilon_0 c^2}$$

or

$$\mathbf{B}_{\text{hole}} = \mathbf{B}_{\text{material}} - \frac{2}{3} \frac{\mathbf{M}}{\epsilon_0 c^2}. \quad (36.35)$$

This result is quite different from what we got for \mathbf{E} .

It is, of course, possible to get these results in a more physical way, by using the Maxwell equations directly. For example, Eq. (36.34) follows directly from $\nabla \cdot \mathbf{B} = 0$. (You use a gaussian surface that is half in the material and half out.) Similarly, you can get Eq. (36.33) by using a line integral along a curve that goes up inside the hole and returns through the material. Physically, the field in the hole is reduced because of the surface currents—which are given by $\nabla \times \mathbf{M}$. We will leave it for you to show that Eq. (36.35) can also be obtained by considering the effects of the surface currents on the boundary of the spherical cavity.

In finding the equilibrium magnetization from Eq. (36.29), it turns out to be most convenient to deal with \mathbf{H} ; so write

$$\mathbf{B}_a = \mathbf{H} + \lambda \frac{\mathbf{M}}{\epsilon_0 c^2}. \quad (36.36)$$

In the spherical hole approximation, we would have $\lambda = \frac{1}{3}$, but, as you will see, we will want later to use some other value, so we leave it as an adjustable parameter. Also, we will take all the fields in the same direction so that we won't need to worry about the vector directions. If we were now to substitute Eq. (36.36) into Eq. (36.29), we would have one equation that relates the magnetization M to the magnetizing field H :

$$M = N\mu \tanh \left(\frac{H + \lambda M / \epsilon_0 c^2}{kT} \right).$$

It is, however, an equation that cannot be solved explicitly, so we will do it graphically.

Let's put the problem in a generalized form by writing Eq. (36.29) as

$$\frac{M}{M_{\text{sat}}} = \tanh x, \quad (36.37)$$

where M_{sat} is the saturation value of the magnetization, namely, $N\mu$, and x represents $\mu B_a / kT$. The dependence of M/M_{sat} on x is shown by curve a in Fig. 36-13. We can also write x as a function of M —using Eq. (36.36) for B_a —as

$$x = \frac{\mu B_a}{kT} = \frac{\mu H}{kT} + \left(\frac{\mu \lambda M_{\text{sat}}}{\epsilon_0 c^2 kT} \right) \frac{M}{M_{\text{sat}}}. \quad (36.38)$$

For any given value of H , this is a straight-line relationship between M/M_{sat} and x . The x intercept is at $x = \mu H / kT$, and the slope is $\epsilon_0 c^2 kT / \mu \lambda M_{\text{sat}}$. For any particular H , we would have a line like the one marked b in Fig. 36-13. The intersection of curves a and b gives us the solution for M/M_{sat} . We have solved the problem.

Let's look at how the solutions will go for various circumstances. We start with $H = 0$. There are two possible situations, shown by the lines b_1 and b_2 in Fig. 36-14. You will notice from Eq. (36.38) that the slope of the line is proportional to the absolute temperature T . So, at *high temperatures* we would have a line like b_1 . The solution is $M/M_{\text{sat}} = 0$. When the magnetizing field H is zero, the magnetization is also zero. But at *low temperatures*, we would have a line like b_2 , and there are *two solutions* for M/M_{sat} —one with $M/M_{\text{sat}} = 0$ and one with M/M_{sat} near one. It turns out that only the upper solution is stable—as you can see by considering small variations about these solutions.

According to these ideas, then, a magnetic material should magnetize itself *spontaneously* at sufficiently low temperatures. In short, when the thermal motions are small enough, the coupling between the atomic magnets causes them all to line up parallel to each other—we have a permanently magnetized material analogous to the ferroelectrics we discussed in Chapter 11.

If we start at high temperatures and come down, there is a critical temperature, called the Curie temperature T_c , where the ferromagnetic behavior suddenly sets in. This temperature corresponds to the line b_3 of Fig. 36-14, which is tangent to the curve a , and has, therefore, a slope of 1. The Curie temperature is given by

$$\frac{\epsilon_0 c^2 k T_c}{\mu \lambda M_{\text{sat}}} = 1. \quad (36.39)$$

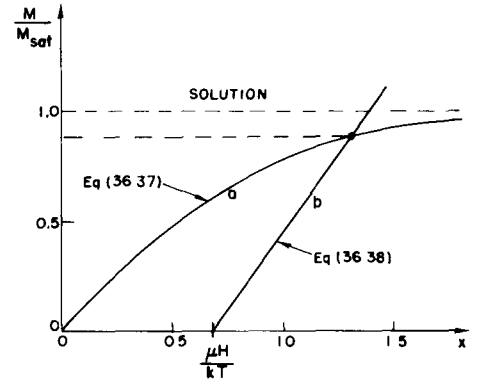


Fig. 36-13. A graphical solution of Eqs. (36.37) and (36.38).

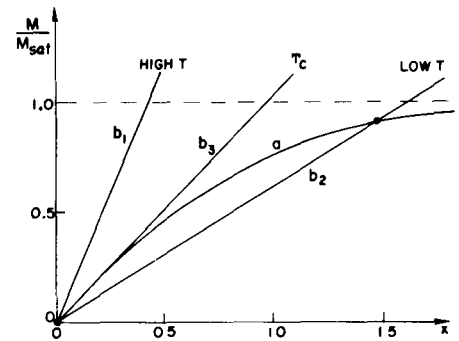


Fig. 36-14. Finding the magnetization when $H = 0$.

We can, if we wish, write Eq. (36.38) more simply in terms of T_c as

$$x = \frac{\mu H}{kT} + \frac{T_c}{T} \left(\frac{M}{M_{\text{sat}}} \right). \quad (36.40)$$

Now we want to see what happens for small magnetizing fields H . We can see from Fig. 36-14 how things will go if we shift our straight lines a little to the right. For the low-temperature case, the intersection point will move out a little bit along the low-slope part of curve a , and M will change relatively little. For the high-temperature case, however, the intersection point runs up the steep part of curve a , and M will change relatively rapidly. In fact, we can approximate this part of curve a by a straight line of unit slope, and write:

$$\frac{M}{M_{\text{sat}}} = x = \frac{\mu H}{kT} + \frac{T_c}{T} \left(\frac{M}{M_{\text{sat}}} \right).$$

Now we can solve for M/M_{sat} :

$$\frac{M}{M_{\text{sat}}} = \frac{\mu H}{k(T - T_c)}. \quad (36.41)$$

We have a law that is something like the one we had for paramagnetism. For paramagnetism, we had

$$\frac{M}{M_{\text{sat}}} = \frac{\mu B}{kT}. \quad (36.42)$$

One difference now is that we have the magnetization in terms of H , which includes some of the effects of the interaction of the atomic magnets, but the main difference is that the magnetization is inversely proportional to the *difference* between T and T_c , instead of to the absolute temperature T , alone. Neglecting the interactions between neighboring atoms corresponds to taking $\lambda = 0$, which from Eq. (36.39) means taking $T_c = 0$. Then the results are just what we had in Chapter 35.

We can check our theoretical picture with the experimental data for nickel. It is observed experimentally that the ferromagnetic behavior of nickel disappears when its temperature is raised above 631°K . We can compare this with T_c calculated from Eq. (36.39). Remembering that $M_{\text{sat}} = \mu N$, we have

$$T_c = \lambda \frac{N_\mu^2}{k\epsilon_0 c^2}.$$

From the density and atomic weight of nickel, we get

$$N = 9.1 \times 10^{28} \text{ m}^{-3}.$$

Taking μ from Eq. (36.28), and setting $\lambda = \frac{1}{3}$, we get

$$T_c = 0.24^\circ\text{K}.$$

There is a discrepancy of a factor of about 2600! Our theory of ferromagnetism fails completely.

We can try to “patch up” the theory as Weiss did by saying that for some unknown reason λ is not one-third, but $(2600) \times \frac{1}{3}$ —or about 900. It turns out that one gets similar values for other ferromagnetic materials like iron. To see what this means, let’s go back to Eq. (36.36). We see that a large λ means that B_a , the local field on the atom, appears to be much, much larger than we would think. In fact, writing $H = B - M/\epsilon_0 c^2$, we have

$$B_a = B + \frac{(\lambda - 1)M}{\epsilon_0 c^2}.$$

According to our original idea—with $\lambda = \frac{1}{3}$ —the local magnetization M *reduces* the effective field B_a by the amount $-\frac{2}{3}M/\epsilon_0$. Even if our model of a spherical hole were not very good, we would still expect *some* reduction. Instead, to explain

the phenomenon of ferromagnetism, we have to imagine that the magnetization of the field *enhances* the local field by some large factor—like one thousand or more. There doesn't seem to be any reasonable way to manufacture such tremendous fields at an atom—nor even fields of the proper sign! Clearly, our “magnetic” theory of ferromagnetism is a dismal failure. We must conclude, then, that ferromagnetism has to do with some *nonmagnetic* interaction between the spinning electrons in neighboring atoms. This interaction must generate a strong tendency for all of the nearby spins to line up in one direction. We will see later that it has to do with quantum mechanics and the Pauli exclusion principle.

Finally, we look at what happens at low temperatures—for $T < T_c$. We have seen that there will then be a spontaneous magnetization—even with $H = 0$ —given by the intersection of the curves a and b_2 of Fig. 36-14. If we solve for M for various temperatures—by varying the slope of the line b_2 —we get the theoretical curve shown in Fig. 36-15. This curve should be the same for all ferromagnetic materials for which the atomic moment comes from a single electron. The curves for other materials are only slightly different.

In the limit, as T goes to absolute zero, M goes to M_{sat} . As the temperature is increased, the magnetization decreases, falling to zero at the Curie temperature. The points in Fig. 36-15 are the experimental observations for nickel. They fit the theoretical curve fairly well. Even though we don't understand the basic mechanism, the general features of the theory seem to be correct.

Finally, there is one more disturbing discrepancy in our attempt to understand ferromagnetism. We have found that above some temperature the material should behave like a paramagnetic substance with a magnetization M proportional to H (or B), and that below that temperature it should become spontaneously magnetized. But that's not what we found when we measured the magnetization curve for iron. It only became permanently magnetized *after* we had “magnetized” it. According to the ideas just discussed, it would magnetize itself! What is wrong? Well, it turns out that if you look at a *small enough crystal* of iron or nickel, it is indeed completely magnetized! But in large pieces of iron, there are many small regions or “domains” that are magnetized in different directions, so that on a large scale the *average* magnetization appears to be zero. In each small domain, however, the iron has a locked-in magnetization with M nearly equal to M_{sat} . The consequences of this domain structure are that gross properties of large pieces of material are quite different from the microscopic properties that we have really been treating. We will take up in the next lecture the story of the practical behavior of bulk magnetic materials.

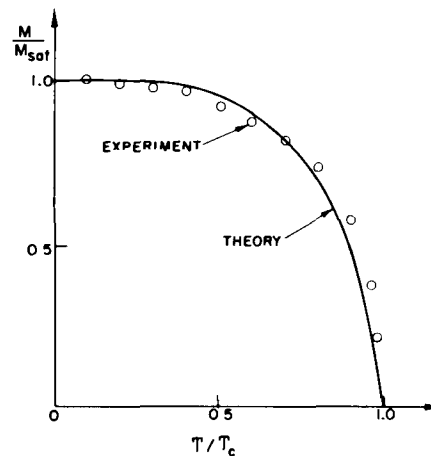


Fig. 36-15. Spontaneous magnetization as a function of temperature for nickel.

Magnetic Materials

37-1 Understanding ferromagnetism

In this chapter we will discuss the behavior and peculiarities of ferromagnetic materials and of other strange magnetic materials. Before proceeding to study magnetic materials, however, we will review very quickly some of the things about the general theory of magnets that we learned in the last chapter.

First, we imagine the atomic currents inside the material that are responsible for the magnetism, and then describe them in terms of a volume current density $\mathbf{j}_{\text{mag}} = \nabla \times \mathbf{M}$. We emphasize that this is not supposed to represent the *actual* currents. When the magnetization is uniform the currents do not *really* cancel out precisely; that is, the whirling currents of one electron in one atom and the whirling currents of an electron in another atom do not overlap in such a way that the sum is exactly zero. Even within a single atom the distribution of magnetism is *not* smooth. For instance, in an iron atom the magnetization is distributed in a more or less spherical shell, not too close to the nucleus and not too far away. Thus, magnetism in matter is quite a complicated thing in its details; it is very irregular. However, we are obliged now to ignore this detailed complexity and discuss phenomena from a gross, average point of view. Then it is true that the *average* current in the interior region, over any finite area that is big compared with an atom, is zero when $\mathbf{M} = 0$. So, what we mean by magnetization per unit volume and \mathbf{j}_{mag} and so on, at the level we are now considering, is an average over regions that are large compared with the space occupied by a single atom.

In the last chapter, we also discovered that a ferromagnetic material has the following interesting property: above a certain temperature it is not strongly magnetic, whereas below this temperature it becomes magnetic. This fact is easily demonstrated. A piece of nickel wire at room temperature is attracted by a magnet. However, if we heat it above its Curie temperature with a gas flame, it becomes nonmagnetic and is not attracted toward the magnet—even when brought quite close to the magnet. If we let it lie near the magnet while it cools off, at the instant its temperature falls below the critical temperature it is suddenly attracted again by the magnet!

The general theory of ferromagnetism that we will use supposes that the spin of the electron is responsible for the magnetization. The electron has spin one-half and carries one Bohr magneton of magnetic moment $\mu = \mu_B = q_e \hbar / 2m$. The electron spin can be pointed either “up” or “down.” Because the electron has a negative charge, when its spin is “up” it has a *negative* moment, and when its spin is “down” it has a *positive* moment. With our usual conventions, the moment μ of the electron is opposite its spin. We have found that the energy of orientation of a magnetic dipole in a given applied field \mathbf{B} is $-\mu \cdot \mathbf{B}$, but the energy of the spinning electrons depends on the neighboring spin alignments as well. In iron, if the moment of a nearby atom is “up,” there is a very strong tendency that the moment of the one next to it will also be “up.” That is what makes iron, cobalt, and nickel so strongly magnetic—the moments all want to be parallel. The first question we have to discuss is *why*.

Soon after the development of quantum mechanics, it was noticed that there is a very strong *apparent* force—not a magnetic force or any other kind of actual force, but only an apparent force—trying to line the spins of nearby electrons *opposite* to one another. These forces are closely related to chemical valence forces. There is a principle in quantum mechanics—called the *exclusion principle*—that

37-1 Understanding ferromagnetism

37-2 Thermodynamic properties

37-3 The hysteresis curve

37-4 Ferromagnetic materials

37-5 Extraordinary magnetic materials

References: Bozorth, R. M., “Magnetism,” *Encyclopaedia Britannica*, Vol. 14, 1957, pp. 636–667.

Kittel, C., *Introduction to Solid State Physics*, John Wiley and Sons, Inc., New York, 2nd ed., 1956.

two electrons cannot occupy exactly the same state, that they cannot be in exactly the same condition as to location and spin orientation.* For example, if they are at the same point, the only alternative is to have their spins opposite. So, if there is a region of space between atoms where electrons like to congregate (as in a chemical bond) and we want to put another electron on top of one already there, the only way to do it is to have the spin of the second one pointed opposite to the spin of the first one. To have the spins parallel is against the law, unless the electrons stay away from each other. This has the effect that a pair of parallel-spin electrons near to each other have much more energy than a pair of opposite-spin electrons; the net effect is as though there were a force trying to turn the spin over. Sometimes this spin-turning force is called the *exchange force*, but that only makes it more mysterious—it is not a very good term. It is just because of the exclusion principle that electrons have a tendency to make their spins opposite. In fact, that is the explanation of the *lack* of magnetism in almost all substances! The spins of the free electrons on the outside of the atoms have tremendous tendency to balance in opposite directions. The problem is to explain why for materials like iron it is just the reverse of what we should expect.

We have summarized the supposed alignment effect by adding a suitable term in the energy equation, by saying that if the electron magnets in the neighborhood have a mean magnetization M , then the moment of an electron has a strong tendency to be in the same direction as the average magnetization of the atoms in the neighborhood. Thus, we may write for the two possible spin orientations,†

$$\begin{aligned}\text{Spin "up" energy} &= +\mu \left(H + \frac{\lambda M}{\epsilon_0 c^2} \right), \\ \text{Spin "down" energy} &= -\mu \left(H + \frac{\lambda M}{\epsilon_0 c^2} \right).\end{aligned}\tag{37.1}$$

When it was clear that quantum mechanics could supply a tremendous spin-orientating force—even if, apparently, of the wrong sign—it was suggested that ferromagnetism might have its origin in this same force, that due to the complexities of iron and the large number of electrons involved, the sign of the interaction energy would come out the other way around. Since the time this was thought of—in about 1927 when quantum mechanics was first being understood—many people have been making various estimates and semicalculations, trying to get a theoretical prediction for λ . The most recent calculations of the energy between the two electron spins in iron—assuming that the interaction is a direct one between the two electrons in neighboring atoms—still give the *wrong sign*. The present understanding of this is again to assume that the complexity of the situation is somehow responsible and to hope that the next man who makes the calculation with a more complicated situation will get the right answer!

It is believed that the up-spin of one of the electrons in the inside shell, which is making the magnetism, tends to make the conduction electrons which fly around the outside have the opposite spin. One might expect this to happen because the conduction electrons come into the same region as the “magnetic” electrons. Since they move around, they can carry their prejudice for being upside down over to the next atom; that is, one “magnetic” electron tries to force the conduction electrons to be opposite, and the conduction electron then makes the next “magnetic” electron opposite to *it*. The double interaction is equivalent to an interaction which tries to line up the two “magnetic” electrons. In other words, the tendency to make parallel spins is the result of an intermediary that tends to some extent to be opposite to both. This mechanism does not require that the conduction electrons be completely “upside down.” They could just have a slight prejudice to be down, just enough to load the “magnetic” odds the other way. This is the mechanism that

* See Chapter 43.

† We write these equations with $H = B - M/\epsilon_0 c^2$ instead of B to agree with the work of the last chapter. You might prefer to write $U = \pm \mu B_a = \pm \mu (B + \lambda' M/\epsilon_0 c^2)$, where $\lambda' = \lambda - 1$. It's the same thing.

the people who have calculated such things now believe is responsible for ferromagnetism. But we must emphasize that to this day nobody can calculate the magnitude of λ simply by knowing that the material is number 26 in the periodic table. In short, we don't thoroughly understand it.

Now let us continue with the theory, and then come back later to discuss a certain error involved in the way we have set it up. If the magnetic moment of a certain electron is "up," energy comes both from the external field and also from the tendency of the spins to be parallel. Since the energy is lower when the spins are parallel, the effect is sometimes thought of as due to an "effective internal field." But remember, it is *not* due to a true *magnetic* force; it is an interaction that is more complicated. In any case, we take Eqs. (37.1) as the formulas for the energies of the two spin states of a "magnetic" electron. At a temperature T , the relative probability of these two states is proportional to $e^{-\text{energy}/kT}$, which we can write as e^{+x} , with $x = \mu(H + \lambda M/\epsilon_0 c^2)/kT$. Then, if we calculate the mean value of the magnetic moment, we find (as in the last chapter) that it is

$$M = N\mu \tanh x. \quad (37.2)$$

Now we would like to calculate the internal energy of the material. We note that the energy of an electron is exactly proportional to the magnetic moment, so that the calculation of the mean moment and the calculation of the mean energy are the same—except that in place of μ in Eq. (37.2) we would write $-\mu B$, which is $-\mu(H + \lambda M/\epsilon_0 c^2)$. The mean energy is then

$$\langle U \rangle_{av} = -N\mu \left(H + \frac{\lambda M}{\epsilon_0 c^2} \right) \tanh x.$$

Now this is not quite correct. The term $\lambda M/\epsilon_0 c^2$ represents interactions of all possible *pairs* of atoms, and we must remember to count each pair only *once*. (When we consider the energy of one electron in the field of the rest and then the energy of a second electron in the field of the rest, we have counted part of the first energy once more.) Thus, we must divide the *mutual interaction term* by two, and our formula for the energy then turns out to be

$$\langle U \rangle_{av} = -N\mu \left(H + \frac{\lambda M}{2\epsilon_0 c^2} \right) \tanh x. \quad (37.3)$$

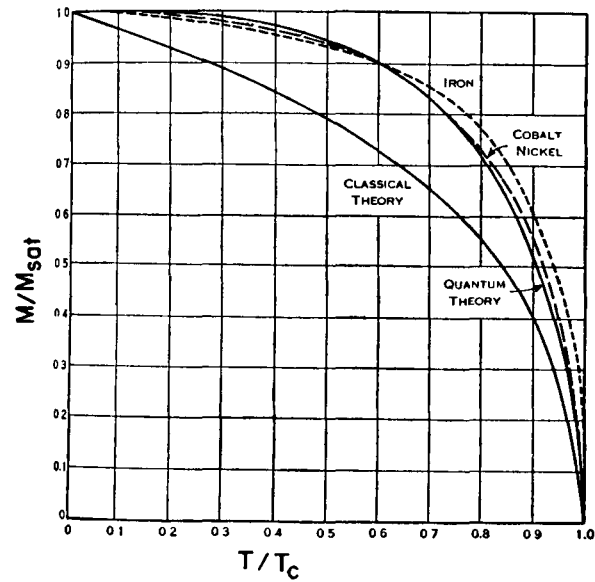
In the last chapter we discovered an interesting thing—that below a certain temperature the material finds a solution to the equations in which the magnetic moment *is not zero*, even with no external magnetizing field. When we set $H = 0$ in Eq. (37.2), we found that

$$\frac{M}{M_{sat}} = \tanh \left(\frac{T_c}{T} \frac{M}{M_{sat}} \right), \quad (37.4)$$

where $M_{sat} = N\mu$, and $T_c = \mu\lambda M_{sat}/k\epsilon_0 c^2$. When we solve this equation (graphically or otherwise), we find that the ratio M/M_{sat} as a function of T/T_c is a curve like that labeled "quantum theory" in Fig. 37-1. The dashed curve marked "cobalt, nickel" shows the experimental results for crystals of these elements. The theory and experiment are in reasonably good agreement. The figure also shows the result of the classical theory in which the calculation is carried out assuming that the atomic magnets can have all possible orientations in space. You can see that this assumption gives a prediction that is not even close to the experimental facts.

Even the quantum theory deviates from the observed behavior at both high and low temperatures. The reason for the deviations is that we have made a rather sloppy approximation in the theory: We have assumed that the energy of an atom depends upon the *mean* magnetization of its neighboring atoms. In other words, for each one that is "up" in the neighborhood of a given atom, there will be a contribution of energy due to that quantum mechanical alignment effect. But how many are there pointed "up"? On the average, that is measured by the

Fig. 37-1. The spontaneous magnetization ($H = 0$) of ferromagnetic crystals as a function of temperature. [Permission from Encyclopaedia Britannica.]



magnetization M —but only *on the average*. A particular atom somewhere might find *all* its neighbors “up.” Then its energy will be larger than the average. Another one might find some up and some down, perhaps averaging to zero, and it would have *no* energy from that term, and so on. What we ought to do is to use some more complicated kind of average, because the atoms in different places have different environments, and the numbers up and down are different for different ones. Instead of just taking one atom subjected to the average influence, we should take each one in its actual situation, compute its energy, and find the *average energy*. But how do we find out how many are “up” and how many are “down” in the neighborhood? That is, of course, just what we are trying to calculate—the number “up” and “down”—so we have a very complicated interconnected problem of correlations, a problem which has never been solved. It is an intriguing and exciting one which has existed for years and on which some of the greatest names in physics have written papers, but even they have not completely solved it.

It turns out that at low temperatures, when almost all the atomic magnets are “up” and only a few are “down,” it is easy to solve; and at high temperatures, far above the Curie temperature T_c when they are almost all random, it is again easy. It is often easy to calculate small departures from some simple, idealized situation, so it is fairly well understood why there are deviations from the simple theory at low temperature. It is also understood physically that for statistical reasons the magnetization *should* deviate at high temperatures. But the exact behavior near the Curie point has never been thoroughly figured out. That’s an interesting problem to work out some day if you want a problem that has never been solved.

37-2 Thermodynamic properties

In the last chapter we laid the groundwork necessary for calculating the thermodynamic properties of ferromagnetic materials. These are, naturally, related to the internal energy of the crystal, which includes interactions of the various spins, given by Eq. (37.3). For the energy of the spontaneous magnetization below the Curie point, we can set $H = 0$ in Eq. (37.3), and—noticing that $\tanh x = M/M_{\text{sat}}$ —we find a mean energy proportional to M^2 :

$$\langle U \rangle_{\text{av}} = - \frac{N\mu\lambda M^2}{2\epsilon_0 c^2 M_{\text{sat}}}. \quad (37.5)$$

If we now plot the energy due to the magnetism as a function of temperature, we get a curve which is the negative of the square of the curve of Fig. 37-1, as drawn in Fig. 37-2(a). If we were to measure then the *specific heat* of such a material we would obtain a curve which is the derivative of 37-2(a). It is shown in Fig.

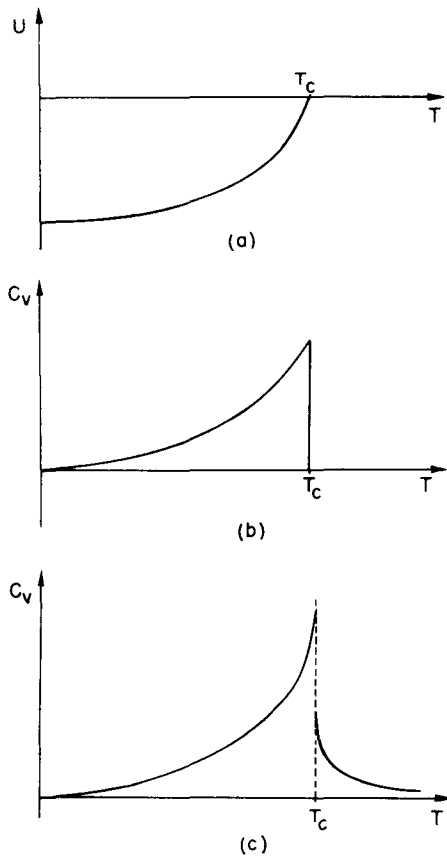


Fig. 37-2. The energy per unit volume and specific heat of a ferromagnetic crystal.

37-2(b). It rises slowly with increasing temperature, but falls suddenly to zero at $T = T_c$. The sharp drop is due to the change in slope of the magnetic energy and is reached right at the Curie point. So without any magnetic measurements at all we could have discovered that something was going on inside of iron or nickel by measuring this thermodynamic property. However, both experiment and improved theory (with fluctuations included) suggest that this simple curve is wrong and that the true situation is really more complicated. The curve goes higher at the peak and falls to zero somewhat slowly. Even if the temperature is high enough to randomize the spins *on the average*, there are still local regions where there is a certain amount of polarization, and in these regions the spins still have a little extra energy of interaction—which only dies out slowly as things get more and more random with further increases in temperature. So the actual curve looks like Fig. 37-2(c). One of the challenges of theoretical physics today is to find an exact theoretical description of the character of the specific heat near the Curie transition—an intriguing problem which has not yet been solved. Naturally, this problem is very closely related to the shape of the magnetization curve in the same region.

Now we want to describe some experiments, other than thermodynamic ones, which show that there is something *right* about our interpretation of magnetism. When the material is magnetized to saturation at low enough temperatures, M is very nearly equal to M_{sat} —nearly all the spins are parallel, as well as their magnetic moments. We can check this by an experiment. Suppose we suspend a bar magnet by a thin fiber and then surround it by a coil so that we can reverse the magnetic field without touching the magnet or putting any torque on it. This is a very difficult experiment because the magnetic forces are so enormous that any irregularities, any lopsidedness, or any lack of perfection in the iron will produce accidental torques. However, the experiment has been done under careful conditions in which such accidental torques are minimized. By means of the magnetic field from a coil that surrounds the bar, we turn all the atomic magnets over at once. When we do this we also change the angular momenta of all the spins from “up” to “down” (see Fig. 37-3). If angular momentum is to be conserved when the spins all turn over, the rest of the bar must have an opposite change in angular momentum. The whole magnet will start to spin. And sure enough, when we do the experiment, we find a slight turning of the magnet. We can measure the total angular momentum given to the whole magnet, and this is simply N times \hbar , the change in the angular momentum of each spin. The ratio of angular momentum to magnetic moment measured this way comes out to within about 10 percent of what we calculate. Actually, our calculations assume that the atomic magnets are due purely to the electron spin, but there is, in addition, some orbital motion also in most materials. The orbital motion is not completely free of the lattice and does not contribute much more than a few percent to the magnetism. As a matter of fact, the saturation magnetic field that one gets taking $M_{sat} = N\mu$ and using the density of iron of 7.9 and the moment μ of the spinning electron is about 20,000 gauss. But according to experiment, it is actually in the neighborhood of 21,500 gauss. This is a typical magnitude of error—5 or 10 percent—due to neglecting the contributions of the orbital moments that have not been included in making the analysis. Thus, a slight discrepancy with the gyromagnetic measurements is quite understandable.

37-3 The hysteresis curve

We have concluded from our theoretical analysis that a ferromagnetic material should spontaneously become magnetized below a certain temperature so that all the magnetism would be in the same direction. But we know that this is not true for an ordinary piece of *unmagnetized* iron. Why isn't all iron magnetized? We can explain it with the help of Fig. 37-4. Suppose the iron were all a big single crystal of the shape shown in Fig. 37-4(a) and spontaneously magnetized all in one direction. Then there would be a considerable external magnetic field, which would have a lot of energy. We can reduce that field energy if we arrange that one side of

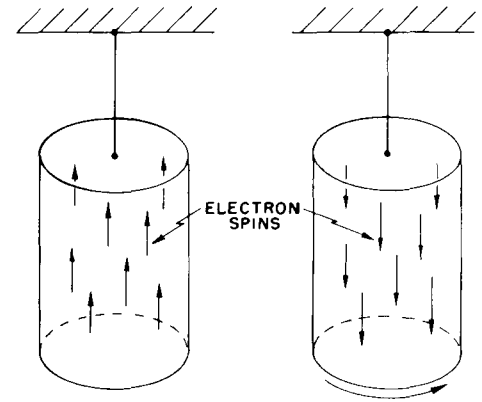


Fig. 37-3. When the magnetization of a bar of iron is reversed, the bar is given some angular velocity.

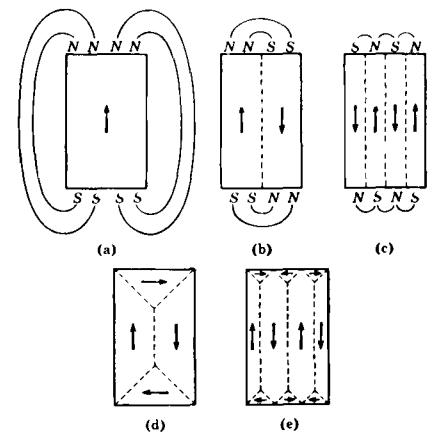


Fig. 37-4. The formation of domains in a single crystal of iron. [From Charles Kittel, *Introduction to Solid State Physics*, John Wiley and Sons, Inc., New York, 2nd ed., 1956.]

the block is magnetized “up” and the other side magnetized “down,” as in Fig. 37-4(b). Then, of course, the fields outside the iron would extend over less volume, so there would be less energy there.

Ah, but wait! In the layer between the two regions we have up-spinning electrons adjacent to down-spinning electrons. But ferromagnetism appears only in those materials for which the energy is *reduced* if the electrons are *parallel* rather than opposite. So, we have added some extra energy along the dotted line in Fig. 37-4(b); this energy is sometimes called *wall energy*. A region having only one direction of magnetization is called a *domain*. At the interface—the “wall”—between two domains, where we have atoms on opposite sides which are spinning in different directions, there is an energy per unit area of the wall. We have described it as though two adjacent atoms were spinning exactly opposite, but it turns out that nature adjusts things so that the transition is more gradual. But we don’t need to worry about such fine details at this point.

Now the question is: When is it better or worse to make a wall? The answer is that it depends on the *size* of the domains. Suppose that we were to scale up a block so that the whole thing was twice as big. The volume in the space outside filled with a given magnetic field strength would be *eight* times bigger, and the energy in the magnetic field, which is proportional to the volume, would also be eight times greater. But the *surface* area between two domains, which will give the wall energy, would be only *four* times as big. Therefore, if the piece of iron is big enough, it will pay to split it into more domains. This is why only the very tiny crystals can have but a single domain. Any large object—one more than about a hundredth of a millimeter in size—will have at least one domain wall; and any ordinary, “centimeter-size” object will be split into many domains, as shown in the figure. Splitting into domains goes on *until the energy needed to put in one extra wall is as large as the energy decrease in the magnetic field outside the crystal*.

Actually nature has discovered still another way to lower the energy: It is not necessary to have the field go outside at all, if a little triangular region is magnetized *sideways*, as in Fig. 37-4(d).^{*} Then with the arrangement of Fig. 37-4(d) we see that there is *no* external field, but instead only a little more domain wall.

But that introduces a new kind of problem. It turns out that when a single crystal of iron is magnetized, it changes its length in the direction of magnetization, so an “ideal” cube with its magnetization, say, “up,” is no longer a perfect cube. The “vertical” dimension will be different from the “horizontal” dimension. This effect is called *magnetostriction*. Because of such geometric changes, the little triangular pieces of Fig. 37-4(d) do not, so to speak, “fit” into the available space anymore—the crystal has got too long one way and too short the other way. Of course, it *does* fit, really, but only by being squashed in; and this involves some mechanical stresses. So, this arrangement *also* introduces an extra energy. It is the balance of all these various energies which determines how the domains finally arrange themselves in their complicated fashion in a piece of unmagnetized iron.

Now, what happens when we put on an external magnetic field? To take a simple case, consider a crystal whose domains are as shown in Fig. 37-4(d). If we apply an external magnetic field in the upward direction, in what manner does the crystal become magnetized? First, the middle domain wall can *move over sideways* (to the right) and reduce the energy. It moves over so that the region which is “up” becomes bigger than the region which is “down”. There are more elementary magnets lined up with the field, and this gives a lower energy. So, for a piece of iron in weak fields—at the very beginning of magnetization—the domain walls begin to move and eat into the regions which are magnetized opposite to the field. As the field continues to increase, a whole crystal shifts gradually into a single

^{*} You may be wondering how spins that have to be either “up” or “down” can also be “sideways”! That’s a good question, but we won’t worry about it right now. We’ll simply adopt the classical point of view, thinking of the atomic magnets as classical dipoles which can be polarized sideways. Quantum mechanics requires considerable expertness to understand how things can be quantized both “up-and-down,” and “right-and-left,” all at the same time.

large domain which the external field helps to keep lined up. In a strong field the crystal “likes” to be all one way *just because* its energy in the applied field is reduced—it is no longer merely the crystal’s own external field which matters.

What if the geometry is not so simple? What if the axes of the crystal and its spontaneous magnetization are in one direction, but we apply the magnetic field in *some other direction*—say at 45° ? We might think that domains would reform themselves with their magnetization parallel to the field, and then as before, they could all grow into one domain. But this is not easy for the iron to do, *for the energy needed to magnetize a crystal depends on the direction of magnetization relative to the crystal axis*. It is relatively easy to magnetize iron in a direction parallel to the crystal axes, but it takes *more* energy to magnetize it in some other direction—like 45° with respect to one of the axes. Therefore, if we apply a magnetic field in such a direction, what happens first is that the domains which point along one of the preferred directions which is *near* to the applied field grow until the magnetization is all along one of these directions. Then *with much stronger fields*, the magnetization is gradually pulled around parallel to the field, as sketched in Fig. 37-5.

In Fig. 37-6 are shown some observations of the magnetization curves of single crystals of iron. To understand them, we must first explain something about the notation that is used in describing directions in a crystal. There are many ways in which a crystal can be sliced so as to produce a face which is a plane of atoms. Everyone who has driven past an orchard or vineyard knows this—it is fascinating to watch. If you look one way, you see lines of trees—if you look another way, you see different lines of trees, and so on. In a similar way, a crystal has definite families of planes that hold many atoms, and the planes have this important characteristic (we consider a cubic crystal to make it easier): If we observe where the planes intersect the three coordinate axes—we find that the *reciprocals* of the three distances from the origin are in the ratio of simple whole numbers. These three whole numbers are taken as the definition of the planes. For example, in Fig. 37-7(a), a plane parallel to the yz -plane is shown. This is called a $[100]$ plane; the reciprocals of its intersection of the y - and z -axes are both zero. The direction perpendicular to such a plane (in a cubic crystal) is given the same set of numbers. It is easy to understand the idea in a cubic crystal, for then the indices $[100]$ mean a vector which has a unit component in the x -direction and none in the y - or z -directions. The $[110]$ direction is in a direction 45° from the x - and y -axes, as in Fig. 37-7(b); and the $[111]$ direction is in the direction of the cube diagonal, as in Fig. 37-7(c).

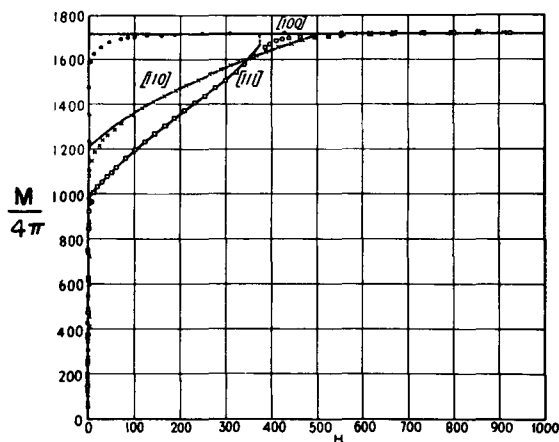


Fig. 37-6. The component of M parallel to H , for different directions of H (with respect to the crystal axes). [From F. Bitter, *Introduction to Ferromagnetism*, McGraw-Hill Book Co., Inc., 1937.]

Returning now to Fig. 37-6, we see the magnetization curves of a single crystal of iron for various directions. First, note that for very tiny fields—so weak that it is hard to see them on the scale at all—the magnetization increases extremely rapidly to quite large values. If the field is in the $[100]$ direction—namely along one of those nice, easy directions of magnetization—the curve goes up to a high value, curves around a little, and then is saturated. What happened is that the

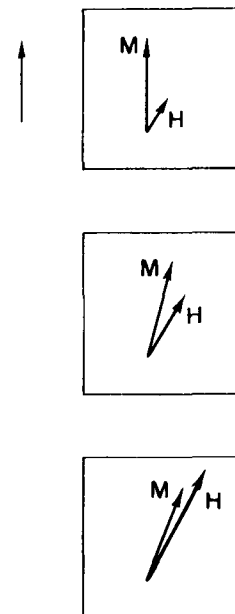


Fig. 37-5. A magnetizing field H at an angle with respect to the crystal axis will gradually change the direction of the magnetization without changing its magnitude.

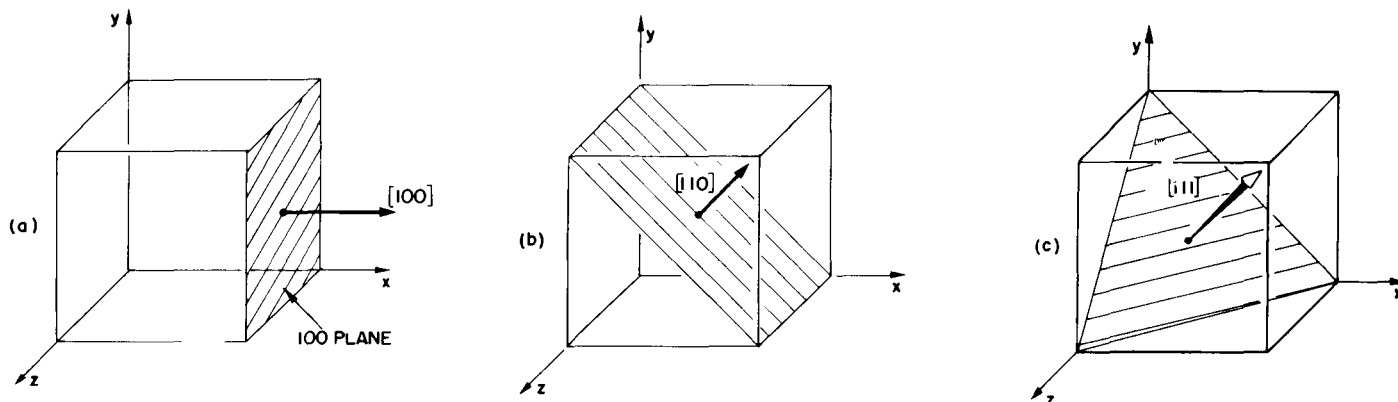


Fig. 37-7. The way the crystal planes are labeled.

domains which were already there are very easily removed. Only a small field is required to make the domain walls move and eat up all of the “wrong-way” domains. Single crystals of iron are enormously permeable (magnetic sense), much more so than ordinary polycrystalline iron. A perfect crystal magnetizes extremely easily. Why is it curved at all? Why doesn’t it just go right up to saturation? We are not sure. You might study that some day. We do understand why it is flat for high fields. When the whole block is a single domain, the extra magnetic field cannot make any more magnetization—it is already at M_{sat} , with all the electrons lines up.

Now, if we try to do the same thing in the $[110]$ direction—which is at 45° to the crystal axes—what will happen? We turn on a little bit of field and the magnetization leaps up as the domains grow. Then as we increase the field some more, we find that it takes quite a lot of field to get up to saturation, because now the magnetization is *turning away* from an “easy” direction. If this explanation is correct, the point at which the $[110]$ curve extrapolates back to the vertical axis should be at $1/\sqrt{2}$ of the saturation value. It turns out, in fact, to be very, very close to $1/\sqrt{2}$. Similarly, in the $[111]$ direction—which is along the cube diagonal—we find, as we would expect, that the curve extrapolates back to nearly $1/\sqrt{3}$ of saturation.

Figure 37-8 shows the corresponding situation for two other materials, nickel and cobalt. Nickel is different from iron. In nickel, it turns out that the $[111]$ direction is the easy direction of magnetization. Cobalt has a hexagonal crystal form, and people have botched up the system of nomenclature for this case. They want to have three axes on the bottom of the hexagon and one perpendicular to these, so they have used four indices. The $[0001]$ direction is the direction of the axis of the hexagon, and $[10\bar{1}0]$ is perpendicular to that axis. We see that crystals of different metals behave in different ways.

Now we must discuss a polycrystalline material, such as an ordinary piece of iron. Inside such materials there are many, many little crystals with their crystal-line axes pointing every which way. *These are not the same as domains.* Remember that the domains were all part of a *single crystal*, but in a piece of iron there are

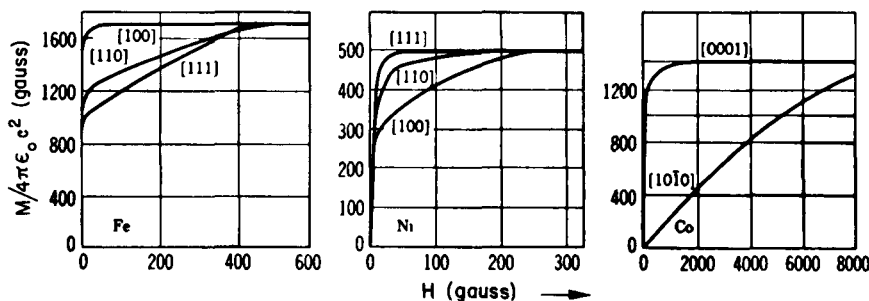


Fig. 37-8. Magnetization curves for single crystals of iron, nickel, and cobalt. [From Charles Kittel, *Introduction to Solid State Physics*, John Wiley and Sons, Inc., New York, 2nd ed., 1956.]

many *different crystals* with axes at different orientations, as shown in Fig. 37-9. Within each of these crystals, there will also generally be some domains. When we apply a *small* magnetic field to a piece of polycrystalline material, what happens is that the domain walls begin to move, and the domains which have a favorable direction of easy magnetization grow larger. This growth is reversible so long as the field stays very small—if we turn the field off, the magnetization will return to zero. This part of the magnetization curve is marked *a* in Fig. 37-10.

For larger fields—in the region *b* of the magnetization curve shown—things get much more complicated. In every small crystal of the material, there are strains and dislocations; there are impurities, dirt, and imperfections. And at all but the smallest fields, the domain wall, in moving, gets stuck on these. There is an interaction energy between the domain wall and a dislocation, or a grain boundary, or an impurity. So when the wall gets to one of them, it gets stuck; it sticks there at a certain field. But then if the field is raised some more, the wall suddenly snaps past. So the motion of the domain wall is not smooth the way it is in a perfect crystal—it gets hung up every once in a while and moves in jerks. If we were to look at the magnetization on a microscopic scale, we would see something like the insert of Fig. 37-10.

Now the important thing is that these jerks in the magnetization can cause an energy loss. In the first place, when a boundary finally slips past an impediment, it moves very quickly to the next one, since the field is already above what would be required for the unimpeded motion. The rapid motion means that there are rapidly changing magnetic fields which produce eddy currents in the crystal. These currents lose energy in heating the metal. A second effect is that when a domain suddenly changes, part of the crystal changes its dimensions from the magnetostriction. Each sudden shift of a domain wall sets up a little sound wave that carries away energy. Because of such effects, the second part of magnetization curve is *irreversible*, and *there is energy being lost*. This is the origin of the hysteresis effect, because to move a boundary wall forward—snap—and then to move it backward—snap—produces a different result. It's like "jerky" friction, and it takes energy.

Eventually, for high enough fields, when we have moved all the domain walls and magnetized each crystal in its best direction, there are still some crystallites which happen to have their easy directions of magnetization not in the direction of our external magnetic field. Then it takes a lot of extra field to turn those magnetic moments around. So the magnetization increases slowly, but smoothly, for high fields—namely in the region marked *c* in the figure. The magnetization does not come sharply to its saturation value, because in the last part of the curve the atomic magnets are *turning* in the strong field. So we see why the magnetization curve of an ordinary polycrystalline materials, such as the one shown in Fig. 37-10, rises a little bit and *reversibly* at first, then rises irreversibly, and then curves over slowly. Of course, there is no sharp break-point between the three regions—they blend smoothly, one into the other.

It is not hard to show that the magnetization process in the middle part of the magnetization curve is jerky—that the domain walls jerk and snap as they shift. All you need is a coil of wire—with many thousands of turns—connected to an amplifier and a loudspeaker, as shown in Fig. 37-11. If you put a few silicon steel sheets (of the type used in transformers) at the center of the coil and bring a bar magnet slowly near the stack, the sudden changes in magnetization will produce impulses of emf in the coil, which are heard as distinct clicks in the loudspeaker. As you move the magnet nearer to the iron you will hear a whole rush of clicks that sound something like the noise of sand grains falling over each other as a can of sand is tilted. The domain walls are jumping, snapping, and jiggling as the field is increased. This phenomenon is called the *Barkhausen effect*.

As you move the magnet even closer to the iron sheets, the noise grows louder and louder for a while but then there is relatively little noise when the magnet gets very close. Why? Because nearly all the domain walls have moved as far as they can go. Any greater field is merely *turning* the magnetization in each domain, which is a smooth process.

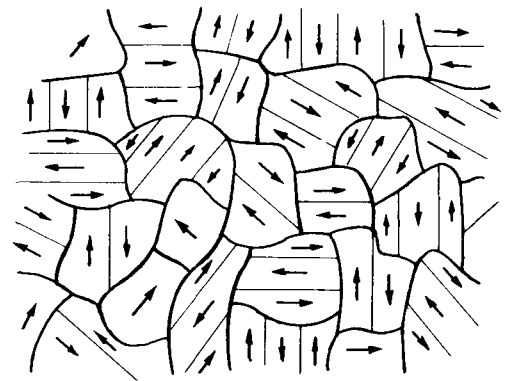


Fig. 37-9. The microscopic structure of an unmagnetized ferromagnetic material. Each crystal grain has an easy direction of magnetization and is broken up into domains which are spontaneously magnetized (usually) parallel to this direction.

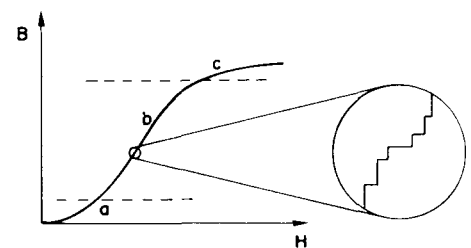


Fig. 37-10. The magnetization curve for polycrystalline iron.

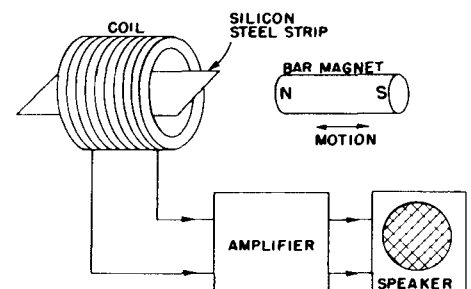


Fig. 37-11. The sudden changes in the magnetization of the steel strip are heard as clicks in the loudspeaker.

If you now withdraw the magnet, so as to come back on the downward branch of the hysteresis loop, the domains all try to get back to low energy again, and you hear another rush of backward-going jerks. You can also note that if you bring the magnet to a given place and move it back and forth a little bit, there is relatively little noise. It is again like tilting a can of sand—once the grains shift into place, small movements of the can don't disturb them. In the iron the small variations in the magnetic field aren't enough to move any boundaries over any of the "humps."

37-4 Ferromagnetic materials

Now we would like to talk about the various kinds of magnetic materials that there are in the technical world and to consider some of the problems involved in designing magnetic materials for different purposes. First, the term "the magnetic properties of iron," which one often hears, is a misnomer—there is no such thing. "Iron" is not a well-defined material—the properties of iron depend critically on the amount of impurities and also on *how* the iron is formed. You can appreciate that the magnetic properties will depend on how easily the domain walls move and that this is a *gross* property, not a property of the individual atoms. So practical ferromagnetism is not really a property of an iron *atom*—it is a property of *solid iron* in a *certain form*. For example, iron can take on two different crystalline forms. The common form has a body-centered cubic lattice, but it can also have a face-centered cubic lattice, which is, however, stable only at temperatures above 1100°C. Of course, at that temperature the body-centered cubic structure is already past the Curie point. However, by alloying chromium and nickel with the iron (one possible mixture is 18 percent chromium and 8 percent nickel) we can get what is called stainless steel, which, although it is mainly iron, retains the face-centered lattice even at low temperatures. Because its crystal structure is different, it has completely different magnetic properties. Most kinds of stainless steel are not magnetic to any appreciable degree, although there are some kinds which are somewhat magnetic—it depends on the composition of the alloy. Even when such an alloy is magnetic, it is not *ferromagnetic* like ordinary iron—even though it is mostly just iron.

We would like now to describe a few of the special materials which have been developed for their particular magnetic properties. First, if we want to make a *permanent* magnet, we would like material with an enormously *wide* hysteresis loop so that, when we turn the current off and come down to zero magnetizing field, the magnetization will remain large. For such materials the domain boundaries should be "frozen" in place as much as possible. One such material is the remarkable alloy "Alnico V" (51% Fe, 8% Al, 14% Ni, 24% Co, 3% Cu). (The rather complex composition of this alloy is indicative of the kind of detailed effort that has gone into making good magnets. What patience it takes to mix five things together and test them until you find the most ideal substance!) When Alnico solidifies, there is a "second phase" which precipitates out, making many tiny grains and very high internal strains. In this material, the domain boundaries have a hard time moving at all. In addition to having a precise composition, Alnico is mechanically "worked" in a way that makes the crystals appear in the form of long grains along the direction in which the magnetization is going to be. Then the magnetization will have a natural tendency to be lined up in these directions and will be held there from the anisotropic effects. Furthermore, the material is even cooled in an external magnetic field when it is manufactured, so that the grains will grow with the right crystal orientation. The hysteresis loop of Alnico V is shown in Fig 37-12. You see that it is about 500 times wider than the hysteresis curve for soft iron that we showed in the last chapter in Fig. 36-8.

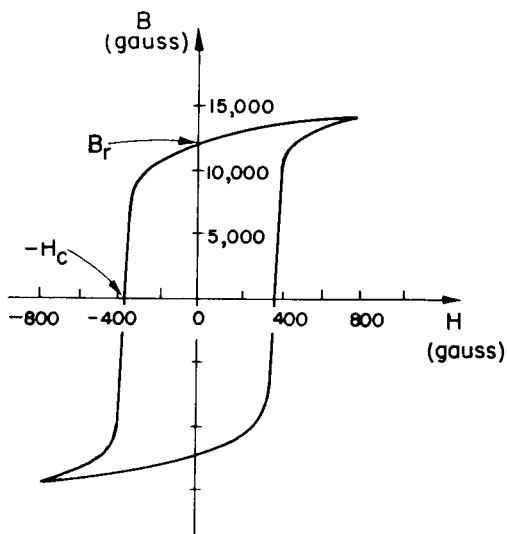


Fig 37-12. The hysteresis curve of Alnico V.

Let's turn now to a different kind of material. For building transformers and motors, we want a material which is magnetically "soft"—one in which the magnetism is easily changed so that an enormous amount of magnetization results from a very small applied field. To arrange this, we need pure, well-annealed material which will have very few dislocations and impurities so that the domain

walls can move easily. It would also be nice if we could make the anisotropy small. Then, even if a grain of the material sits at the wrong angle with respect to the field, it will still magnetize easily. Now we have said that iron prefers to magnetize along the [100] direction, whereas nickel prefers the [111] direction; so if we mix iron and nickel in various proportions, we might hope to find that with just the right proportions the alloy wouldn't prefer *any* direction—the [100] and [111] directions would be equivalent. It turns out that this happens with a mixture of 70 percent nickel and 30 percent iron. In addition—possibly by luck or maybe because of some physical relationship between the anisotropy and the magnetostriction effects—it turns out that the *magnetostriction* of iron and nickel has the opposite sign. And in an alloy of the two metals, this property goes through zero at about 80 percent nickel. So somewhere between 70 and 80 percent nickel we get very “soft” magnetic materials—alloys that are very easy to magnetize. They are called the *permalloys*. Permalloys are useful for high-quality transformers (at low signal levels), but they would be no good at all for permanent magnets. Permalloys must be very carefully made and handled. The magnetic properties of a piece of permalloy are drastically changed if it is stressed beyond its elastic limit—it mustn't be bent. Then, its permeability is reduced because of the dislocations, slip bands, and so on, which are produced by the mechanical deformations. The domain boundaries are no longer easy to move. The high permeability can, however, be restored by annealing at high temperatures.

It is often convenient to have some numbers to characterize the various magnetic materials. Two useful numbers are the intercepts of the hysteresis loop with the B - and H -axes, as indicated in Fig. 37-12. These intercepts are called the *remanent magnetic field* B_r and the *coercive force* H_c . In Table 37-1 we list these numbers for a few magnetic materials.

Table 37-1
Properties of some ferromagnetic materials

Material	B_r Residual magnetic field (gauss)	H_c Coercive force (gauss)
Supermalloy	(≈ 5000)	0.004
Silicon steel (transformer)	12,000	0.05
Armco iron	4000	0.6
Alnico V	13,000	550.

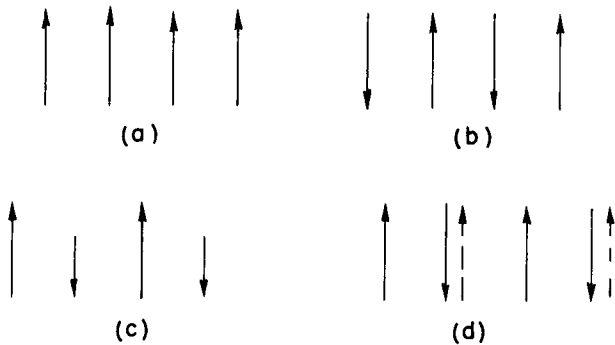


Fig. 37-13. Relative orientation of electron spins in various materials: (a) ferromagnetic, (b) antiferromagnetic, (c) ferrite, (d) yttrium-iron alloy. (Broken arrows show direction of total angular momentum, including orbital motion.)

37-5 Extraordinary magnetic materials

We would now like to discuss some of the more exotic magnetic materials. There are many elements in the periodic table which have incomplete inner electron shells and hence have atomic magnetic moments. For instance, right next to the ferromagnetic elements iron, nickel, and cobalt you will find chromium and manganese. Why aren't *they* ferromagnetic? The answer is that the λ term in Eq. (37.1) has the *opposite sign* for these elements. In the chromium lattice, for example, the spins of the chromium atoms alternate *atom by atom*, as shown in Fig. 37-13(b). So chromium *is* “magnetic” from its own point of view, but it is not technically interesting because there are no *external* magnetic effects. Chromium, then, is an example of a material in which quantum mechanical effects make the spins alternate. Such a material is called *antiferromagnetic*. The alignment in antiferromagnetic materials is also temperature dependent. Below a critical temperature, all the spins are lined up in the alternating array, but when the material is heated above a certain temperature—which is again called the Curie temperature—the spins suddenly become random. There is, internally, a sudden transition. This transition can be seen in the specific heat curve. Also it shows up in some special “magnetic” effects. For instance, the existence of the alternating spins can be verified by scattering neutrons from a crystal of chromium. Because a neutron itself has a spin

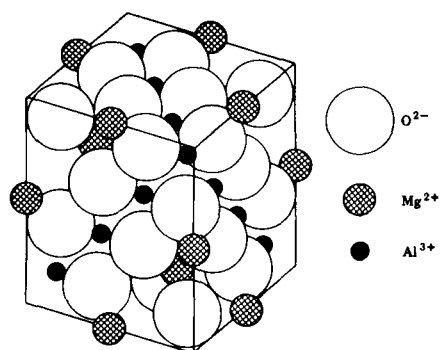


Fig. 37-14. Crystal structure of the mineral spinel (MgAl_2O_4); the Mg^{+2} ions occupy tetrahedral sites, each surrounded by four oxygen ions; the Al^{+3} ions occupy octahedral sites, each surrounded by six oxygen ions. [From Charles Kittel, *Introduction to Solid State Physics*, John Wiley and Sons, Inc., New York, 2nd ed., 1956]

(and a magnetic moment), it has a different amplitude to be scattered, depending on whether its spin is parallel or opposite to the spin of the scatterer. Thus, we get a different interference pattern when the spins in a crystal are alternating than we do when they have a random distribution

There is another kind of substance in which quantum mechanical effects make the electron spins alternate, but which is nevertheless *ferromagnetic*—that is, the crystal has a permanent net magnetization. The idea behind such materials is shown in Fig. 37-14. The figure shows the crystal structure of *spinel*, a magnesium-aluminum oxide, which—as it is shown—is *not* magnetic. The oxide has two kinds of metal atoms: magnesium and aluminum. Now if we replace the magnesium and the aluminum by two magnetic elements like iron and zinc, or by zinc and manganese—in other words, if we put in *magnetic* atoms instead of the nonmagnetic ones—an interesting thing happens. Let's call one kind of metal atom *a* and the other kind of metal atom *b*; then the following combination of forces must be considered. There is an *a-b* interaction which tries to make the *a* atoms and the *b* atoms have opposite spins—because quantum mechanics always gives the opposite sign (except for the mysterious crystals of iron, nickel, and cobalt). Then, there is a direct *a-a* interaction which tries to make the *a*'s opposite, and also a *b-b* interaction which tries to make the *b*'s opposite. Now, of course we cannot have everything opposite everything else—*a* opposite *b*, *a* opposite *a*, and *b* opposite *b*. Presumably because of the distances between the *a*'s and the presence of the oxygen (although we really don't know why), it turns out that the *a-b* interaction is stronger than the *a-a* or the *b-b*. So the solution that nature uses in this case is to make all the *a*'s *parallel* to each other, and all the *b*'s *parallel* to each other, but the two systems *opposite*. That gives the lowest energy because of the stronger *a-b* interaction. The result: all the *a*'s are spinning up and all the *b*'s are spinning down—or vice versa, of course. But if the *magnetic moments* of the *a*-type atom and the *b*-type atom are *not equal*, we can get the situation shown in Fig. 37-13(c), and there can be a net magnetization in the material. The material will then be ferromagnetic—although somewhat weak. Such materials are called *ferrites*. They do not have as high a saturation magnetization as iron—for obvious reasons—so they are only useful for smaller fields. But they have a very important difference—they are insulators; the ferrites are *ferromagnetic insulators*. In high-frequency fields, they will have very small eddy currents and so can be used, for example, in microwave systems. The microwave fields will be able to get inside such an insulating material, whereas they would be kept out by the eddy currents in a conductor like iron.

There is another class of magnetic materials which has only recently been discovered—members of the family of the orthosilicates called *garnets*. They are again crystals in which the lattice contains two kinds of metallic atoms, and we have again a situation in which two kinds of atoms can be substituted almost at will. Among the many compounds of interest there is one which is completely ferromagnetic. It has yttrium and iron in the garnet structure, and the reason it is ferromagnetic is very curious. Here again quantum mechanics is making the neighboring spins opposite, so that there is a locked-in system of spins with the electron spins of the iron one way and the electron spins of the yttrium the opposite way. But the yttrium atom is complicated. It is a rare-earth element and gets a large contribution to its magnetic moment from *orbital* motion of the electrons. For yttrium, the orbital motion contribution is *opposite* that of the spin and also is bigger. Thus, although quantum mechanics, working through the exclusion principle, makes the *spins* of the yttrium opposite those of the iron, it makes the *total* magnetic moment of the yttrium atom *parallel* to the iron because of the orbital effect—as sketched in Fig. 37-13(d). The compound is therefore a regular ferromagnet.

Another interesting example of ferromagnetism occurs in some of the rare-earth elements. It has to do with a still more peculiar arrangement of the spins. The material is not ferromagnetic in the sense that the spins are all parallel, nor is it antiferromagnetic in the sense that every atom is opposite. In these crystals all of the spins *in one layer* are parallel and lie in the plane of the layer. In the next

layer all spins are again parallel to each other, but point in a somewhat different direction. In the following layer they are in still another direction, and so on. The result is that the local magnetization vector varies in the form of a spiral—the magnetic moments of the successive layers rotate as we proceed along a line perpendicular to the layers. It is interesting to try to analyze what happens when a field is applied to such a spiral—all the twistings and turnings that must go on in all those atomic magnets. (Some people *like* to amuse themselves with the theory of these things!) Not only are there cases of “flat” spirals, but there are also cases in which the directions the magnetic moments of successive layers map out a cone, so that it has a spiral component and also a uniform ferromagnetic component in one direction!

The magnetic properties of materials, worked out on a more advanced level than we have been able to do here, have fascinated physicists of all kinds. In the first place, there are those practical people who love to work out ways of making things in a better way—they love to design better and more interesting magnetic materials. The discovery of things like ferrites, or their application, immediately delights people who like to see clever new ways of doing things. Besides this, there are those who find a fascination in the terrible complexity that nature can produce using a few basic laws. Starting with one and the same general idea, nature goes from the ferromagnetism of iron and its domains, to the antiferromagnetism of chromium, to the magnetism of ferrites and garnets, to the spiral structure of the rare earth elements, and on, and on. It is fascinating to discover experimentally all the strange things that go on in these special substances. Then, to the theoretical physicists, ferromagnetism presents a number of very interesting, unsolved, and beautiful challenges. One challenge is to understand why it exists at all. Another is to predict the statistics of the interacting spins in an ideal lattice. Even neglecting any possible extraneous complications, this problem has, so far, defied full understanding. The reason that it is so interesting is that it is such an easily stated problem: Given a lot of electron spins in a regular lattice, interacting with such-and-such a law, what do they do? It is simply stated, but it has defied complete analysis for years. Although it has been analyzed rather carefully for temperatures not too close to the Curie point, the theory of the sudden transition ~~at~~ the Curie point still needs to be completed.

Finally, the whole subject of the system of spinning atomic magnets—in ferromagnetic, or in paramagnetic materials and in nuclear magnetism, has also been a fascinating thing to advanced students in physics. The system of spins can be pushed on and pulled on with external magnetic fields, so one can do many tricks with resonances, with relaxation effects, with spin-echoes, and with other effects. It serves as a prototype of many complicated thermodynamic systems. But in paramagnetic materials the situation is often fairly simple, and people have been delighted both to do experiments and to explain the phenomena theoretically.

We now close our study of electricity and magnetism. In the first chapter, we spoke of the great strides that have been made since the early Greek observation of the strange behaviors of amber and of lodestone. Yet in all our long and involved discussion we have never explained *why it is that when we rub a piece of amber we get a charge on it*, nor have we explained *why a lodestone is magnetized*! You may say, “Oh, we just didn’t get the right sign.” No, it is worse than that. Even if we *did* get the right sign, we would still have the question: Why is the piece of lodestone in the ground magnetized? There is the earth’s magnetic field, of course, but *where does the earth’s field come from*? Nobody really knows—there have only been some good guesses. So you see, this physics of ours is a lot of fakery—we start out with the phenomena of lodestone and amber, and we end up not understanding either of them very well. But we *have* learned a tremendous amount of very exciting and very practical information in the process!

Elasticity

38-1 Hooke's law

The subject of elasticity deals with the behavior of those substances which have the property of recovering their size and shape when the forces producing deformations are removed. We find this elastic property to some extent in all solid bodies. If we had the time to deal with the subject at length, we would want to look into many things: the behavior of materials, the general laws of elasticity, the general theory of elasticity, the atomic machinery that determine the elastic properties, and finally the limitations of elastic laws when the forces become so great that plastic flow and fracture occur. It would take more time than we have to cover all these subjects in detail, so we will have to leave out some things. For example, we will not discuss plasticity or the limitations of the elastic laws. (We touched on these subjects briefly when we were talking about dislocations in metals.) Also, we will not be able to discuss the internal mechanisms of elasticity—so our treatment will not have the completeness we have tried to achieve in the earlier chapters. Our aim is mainly to give you an acquaintance with some of the ways of dealing with such practical problems as the bending of beams.

When you push on a piece of material, it “gives”—the material is deformed. If the force is small enough, the relative displacements of the various points in the material are proportional to the force—we say the behavior is *elastic*. We will discuss only the elastic behavior. First, we will write down the fundamental laws of elasticity, and then we will apply them to a number of different situations.

Suppose we take a rectangular block of material of length l , width w , and height h , as shown in Fig. 38-1. If we pull on the ends with a force F , then the length increases by an amount Δl . We will suppose in all cases that the change in length is a small fraction of the original length. As a matter of fact, for materials like wood and steel, the material will break if the change in length is more than a few percent of the original length. For a large number of materials, experiments show that for sufficiently small extensions the force is proportional to the extension

$$F \propto \Delta l. \quad (38.1)$$

This relation is known as *Hooke's law*.

The lengthening Δl of the bar will also depend on its length. We can figure out how by the following argument. If we cement two identical blocks together, end to end, the same forces act on each block, each will stretch by Δl . Thus, the stretch of a block of length $2l$ would be twice as big as a block of the same cross section, but of length l . In order to get a number more characteristic of the material, and less of any particular shape, we choose to deal with the ratio $\Delta l/l$ of the extension to the original length. This ratio is proportional to the force but independent of l :

$$F \propto \frac{\Delta l}{l}. \quad (38.2)$$

The force F will also depend on the area of the block. Suppose that we put two blocks side by side. Then for a given stretch Δl we would have the force F on each block, or twice as much on the combination of the two blocks. The force, for a given amount of stretch, must be proportional to the cross-sectional area A of the block. To obtain a law in which the coefficient of proportionality is independent of the dimensions of the body, we write Hooke's law for a rectangular

38-1 Hooke's law

38-2 Uniform strains

38-3 The torsion bar; shear waves

38-4 The bent beam

38-5 Buckling

Review: Chapter 47, Vol. I, *Sound; the Wave Equation*.

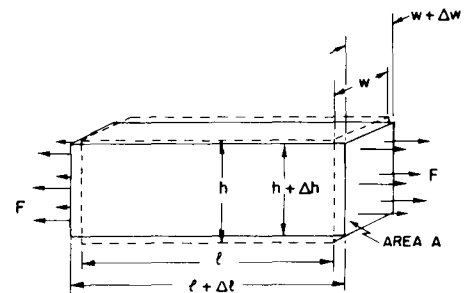


Fig. 38-1. The stretching of a bar under uniform tension.

block in the form

$$F = YA \frac{\Delta l}{l}. \quad (38.3)$$

The constant Y is a property only of the nature of the material; it is known as *Young's modulus*. (Usually you will see Young's modulus called E . But we've used E for electric fields, energy, and emf's, so we prefer to use a different letter.)

The *force per unit area* is called the *stress*, and the stretch per unit length—the *fractional stretch*—is called the *strain*. Equation (38.3) can therefore be rewritten in the following way:

$$\frac{F}{A} = Y \times \frac{\Delta l}{l}, \quad (38.4)$$

$$\text{Stress} = (\text{Young's modulus}) \times (\text{Strain}).$$

There is another part to Hooke's law: When you *stretch* a block of material in one direction it *contracts* at right angles to the stretch. The contraction in width is proportional to the width w and also to $\Delta l/l$. The sideways contraction is in the same proportion for both width and height, and is usually written

$$\frac{\Delta w}{w} = \frac{\Delta h}{h} = -\sigma \frac{\Delta l}{l}, \quad (38.5)$$

where the constant σ is another property of the material called *Poisson's ratio*. It is always positive in sign and is a number less than $1/2$. (It is "reasonable" that σ should be generally positive, but it is not quite clear that it *must* be so.)

The two constants Y and σ specify completely the elastic properties of a *homogeneous isotropic* (that is, noncrystalline) material. In crystalline materials the stretches and contractions can be different in different directions, so there can be many more elastic constants. We will restrict our discussion temporarily to homogeneous isotropic materials whose properties can be described by Y and σ . As usual there are different ways of describing things—some people like to describe the elastic properties of materials by different constants. It always takes two, and they can be related to σ and Y .

The last general law we need is the principle of superposition. Since the two laws (38.4) and (38.5) are linear in the forces and in the displacements, superposition will work. If you have one set of forces and get some displacements, and then you add a new set of forces and get some additional displacements, the resulting displacements will be the sum of the ones you would get with the two sets of forces acting independently.

Now we have all the general principles—the superposition principle and Eqs. (38.4) and (38.5)—and that's all there is to elasticity. But that is like saying that once you have Newton's laws that's all there is to mechanics. Or, given Maxwell's equations, that's all there is to electricity. It is, of course, true that with these principles you have a great deal, because with your present mathematical ability you could go a long way. We will, however, work out a few special applications.

38-2 Uniform strains

As our first example let's find out what happens to a rectangular block under uniform hydrostatic pressure. Let's put a block under water in a pressure tank. Then there will be a force acting inward on every face of the block proportional to the area (see Fig. 38-2). Since the hydrostatic pressure is uniform, the *stress* (force per unit area) on each face of the block is the same. We will work out first the change in the length. The change in length of the block can be thought of as the sum of changes in length that would occur in the three independent problems which are sketched in Fig. 38-3.

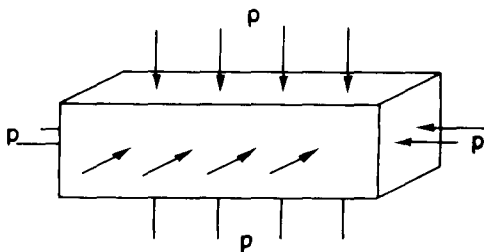


Fig. 38-2. A bar under uniform hydrostatic pressure.

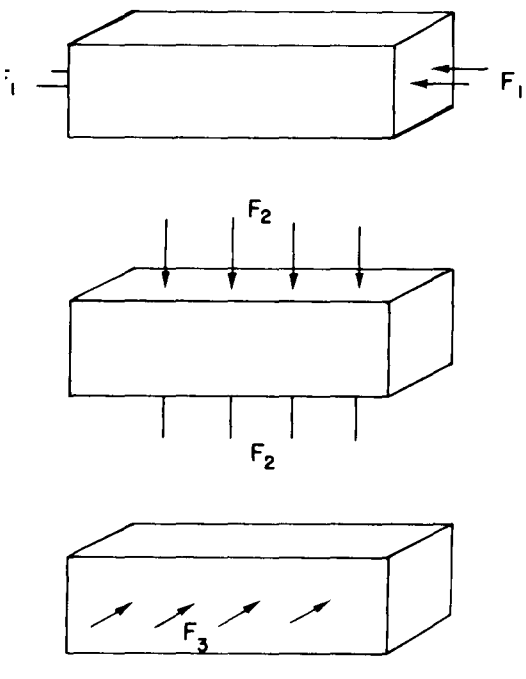


Fig. 38-3. Hydrostatic pressure is the superposition of three longitudinal compressions.

Problem 1. If we push on the ends of the block with a pressure p , the compressional strain is p/Y , and it is negative,

$$\frac{\Delta l_1}{l} = -\frac{p}{Y}.$$

Problem 2. If we push on the two sides of the block with pressure p , the compressional strain is again p/Y , but now we want the lengthwise strain. We can get that from the sideways strain multiplied by $-\sigma$. The sideways strain is

$$\frac{\Delta w}{w} = -\frac{p}{Y};$$

so

$$\frac{\Delta l_2}{l} = +\sigma \frac{p}{Y}.$$

Problem 3. If we push on the top of the block, the compressional strain is once more p/Y , and the corresponding strain in the sideways direction is again $-\sigma p/Y$. We get

$$\frac{\Delta l_3}{l} = +\sigma \frac{p}{Y}.$$

Combining the results of the three problems—that is, taking $\Delta l = \Delta l_1 + \Delta l_2 + \Delta l_3$ —we get

$$\frac{\Delta l}{l} = -\frac{p}{Y} (1 - 2\sigma). \quad (38.6)$$

The problem is, of course, symmetrical in all three directions; it follows that

$$\frac{\Delta w}{w} = \frac{\Delta h}{h} = -\frac{p}{Y} (1 - 2\sigma). \quad (38.7)$$

The change in the *volume* under hydrostatic pressure is also of some interest. Since $V = lwh$, we can write, for small displacements,

$$\frac{\Delta V}{V} = \frac{\Delta l}{l} + \frac{\Delta w}{w} + \frac{\Delta h}{h}.$$

Using (38.6) and (38.7), we have

$$\frac{\Delta V}{V} = -3 \frac{p}{Y} (1 - 2\sigma). \quad (38.8)$$

People like to call $\Delta V/V$ the *volume strain* and write

$$p = -K \frac{\Delta V}{V}.$$

The *volume stress* p is proportional to the volume strain—Hooke's law once more. The coefficient K is called the *bulk modulus*; it is related to the other constants by

$$K = \frac{Y}{3(1 - 2\sigma)}. \quad (38.9)$$

Since K is of some practical interest, many handbooks give Y and K instead of Y and σ . If you want σ you can always get it from Eq. (38.9). We can also see from Eq. (38.9) that Poisson's ratio, σ , must be less than one-half. If it were not, the bulk modulus K would be negative, and the material would expand under increasing pressure. That would allow us to get mechanical energy *out of* any old block—it would mean that the block was in unstable equilibrium. If it started to expand it would continue by itself with a release of energy.

Now we want to consider what happens when you put a “shear” strain on something. By shear strain we mean the kind of distortion shown in Fig. 38-4. As a preliminary to this, let us look at the strains in a *cube* of material subjected to the forces shown in Fig. 38-5. Again we can break it up into two problems: the vertical

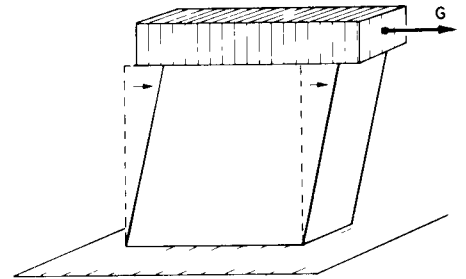


Fig. 38-4 A cube in uniform shear.

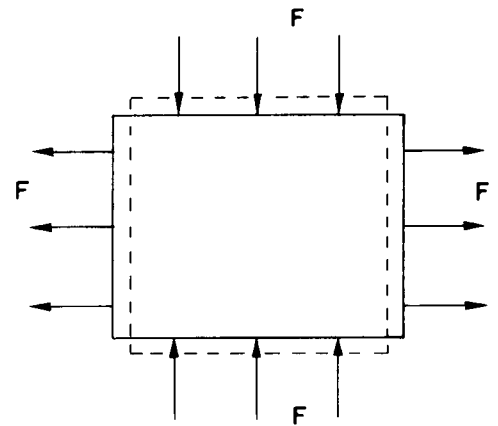


Fig. 38-5. A cube with compressing forces on top and bottom and equal stretching forces on two sides.

pushes, and the horizontal pulls. Calling A the area of the cube face, we have for the change in horizontal length

$$\frac{\Delta l}{l} = \frac{1}{Y} \frac{F}{A} + \sigma \frac{1}{Y} \frac{F}{A} = \frac{1 + \sigma}{Y} \frac{F}{A}. \quad (38.10)$$

The change in the vertical height is just the negative of this.

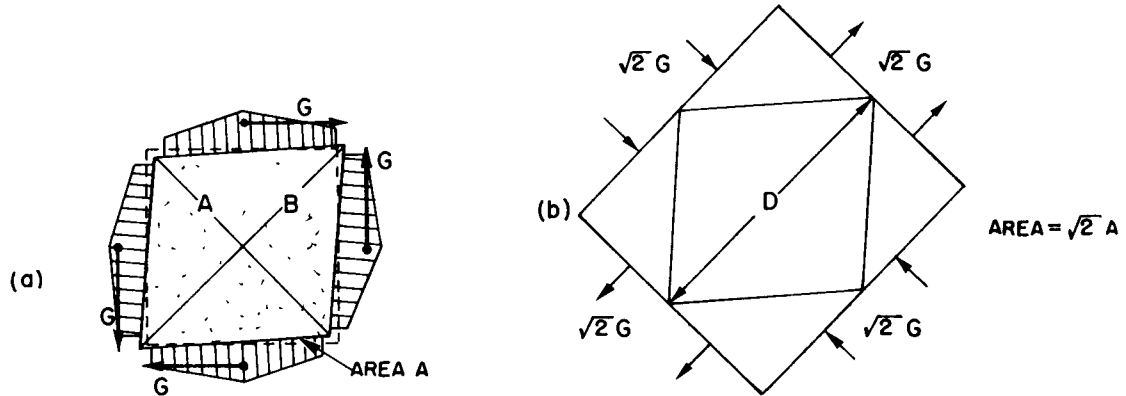


Fig. 38-6. The two pairs of shear forces in (a) produce the same stress as the compressing and stretching forces of (b).

Now suppose we have the same cube and subject it to the shearing forces shown in Fig. 38-6(a). Note that all the forces have to be equal if there are to be no net torques and the cube is to be in equilibrium. (Similar forces must also exist in Fig. 38-4, since the block is in equilibrium. They are provided through the “glue” that holds the block to the table.) The cube is then said to be in a state of pure shear. But note that if we cut the cube by a plane at 45° —say along the diagonal A in the figure—the total force acting across the plane is *normal* to plane and is equal to $\sqrt{2}G$. The area over which this force acts is $\sqrt{2}A$; therefore, the tensile stress normal to this plane is simply G/A . Similarly, if we examine a plane at an angle of 45° the other way—the diagonal B in the figure—we see that there is a compressional stress normal to this plane of $-G/A$. From this, we see that the *stress* in a “pure shear” is equivalent to a combination of tension and compression stresses of equal strength and at right angles to each other, and at 45° to the original faces of the cube. The internal stresses and strains are the same as we would find in the larger block of material with the forces shown in Fig. 38-6(b). But this is the problem we have already solved. The change in length of the diagonal is given by Eq. (38.10),

$$\frac{\Delta D}{D} = \frac{1 + \sigma}{Y} \frac{G}{A}. \quad (38.11)$$

(One diagonal is shortened; the other is elongated.)

It is often convenient to express a shear strain in terms of the angle by which the cube is twisted—the angle θ in Fig. 38-7. From the geometry of the figure you can see that the horizontal shift δ of the top edge is equal to $\sqrt{2} \Delta D$. So

$$\theta = \frac{\delta}{l} = \frac{\sqrt{2} \Delta D}{l} = 2 \frac{\Delta D}{D}. \quad (38.12)$$

The shear stress g is defined as the tangential force on one face divided by the area, $g = G/A$. Using Eq. (38.11) in (38.12), we get

$$\theta = 2 \frac{1 + \sigma}{Y} g.$$

Or, writing this in the form “stress = constant times strain,”

$$g = \mu \theta. \quad (38.13)$$

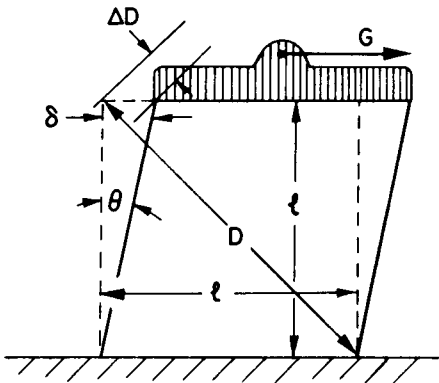


Fig. 38-7. The shear strain θ is $2 \Delta D/D$.

The proportionality coefficient μ is called the *shear modulus* (or, sometimes, the coefficient of rigidity). It is given in terms of Y and σ by

$$\mu = \frac{Y}{2(1 + \sigma)}. \quad (38.14)$$

Incidentally, the shear modulus must be positive—otherwise you could get work out of a self-shearing block. From Eq. (38.14), σ must be greater than -1 . We know, then, that σ must be between -1 and $+\frac{1}{2}$; in practice, however, it is always greater than zero.

As a last example of the type of situation where the stresses are uniform through the material, let's consider the problem of a block which is stretched, while it is at the same time *constrained* so that no lateral contraction can take place. (Technically, it's a little easier to compress it while keeping the sides from bulging out—but it's the same problem.) What happens? Well, there must be sideways forces which keep it from changing its thickness—forces we don't know off-hand but will have to calculate. It's the same kind of problem we have already done, only with a little different algebra. We imagine forces on all three sides, as shown in Fig. 38-8; we calculate the changes in dimensions, and we choose the transverse forces to make the width and height remain constant. Following the usual arguments, we get for the three strains:

$$\frac{\Delta l_x}{l_x} = \frac{1}{Y} \frac{F_x}{A_x} - \frac{\sigma}{Y} \frac{F_y}{A_y} - \frac{\sigma}{Y} \frac{F_z}{A_z} = \frac{1}{Y} \left[\frac{F_x}{A_x} - \sigma \left(\frac{F_y}{A_y} + \frac{F_z}{A_z} \right) \right], \quad (38.15)$$

$$\frac{\Delta l_y}{l_y} = \frac{1}{Y} \left[\frac{F_y}{A_y} - \sigma \left(\frac{F_x}{A_x} + \frac{F_z}{A_z} \right) \right], \quad (38.16)$$

$$\frac{\Delta l_z}{l_z} = \frac{1}{Y} \left[\frac{F_z}{A_z} - \sigma \left(\frac{F_x}{A_x} + \frac{F_y}{A_y} \right) \right]. \quad (38.17)$$

Now since Δl_y and Δl_z are supposed to be zero, Eqs. (38.16) and (38.17) give two equations relating F_y and F_z to F_x . Solving them together, we get that

$$\frac{F_y}{A_y} = \frac{F_z}{A_z} = \frac{\sigma}{1 - \sigma} \frac{F_x}{A_x}. \quad (38.18)$$

Substituting in (38.15), we have

$$\frac{\Delta l_x}{l_x} = \frac{1}{Y} \left(1 - \frac{2\sigma^2}{1 - \sigma} \right) \frac{F_x}{A_x} = \frac{1}{Y} \left(\frac{1 - \sigma - 2\sigma^2}{1 - \sigma} \right) \frac{F_x}{A_x}. \quad (38.19)$$

Often, you will see this turned around, and with the quadratic in σ factored out, it is then written

$$\frac{F}{A} = \frac{1 - \sigma}{(1 + \sigma)(1 - 2\sigma)} Y \frac{\Delta l}{l}. \quad (38.20)$$

When we constrain the sides, Young's modulus gets multiplied by a complicated function of σ . As you can most easily see from Eq. (38.19), the factor in front of Y is always greater than 1. It is harder to stretch the block when the sides are held—which also means that a block is *stronger* when the sides are held than when they are not.

38-3 The torsion bar; shear waves

Let's now turn our attention to an example which is more complicated because different parts of the material are stressed by different amounts. We consider a twisted rod such as you would find in a drive shaft of some machinery, or in a quartz fiber suspension used in a delicate instrument. As you probably know from experiments with the torsion pendulum, the *torque* on a twisted rod is proportional to the *angle*—the constant of proportionality obviously depending upon the length of the rod, on the radius of the rod, and on the properties of the material. The question is: In what way? We are now in a position to answer this question; it's just a matter of working out some geometry.

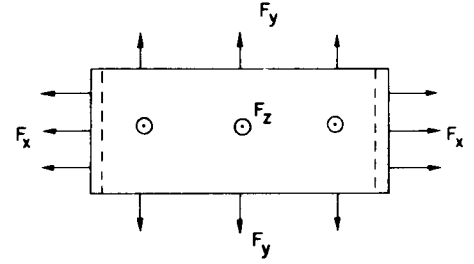


Fig. 38-8. Stretching without lateral contraction.

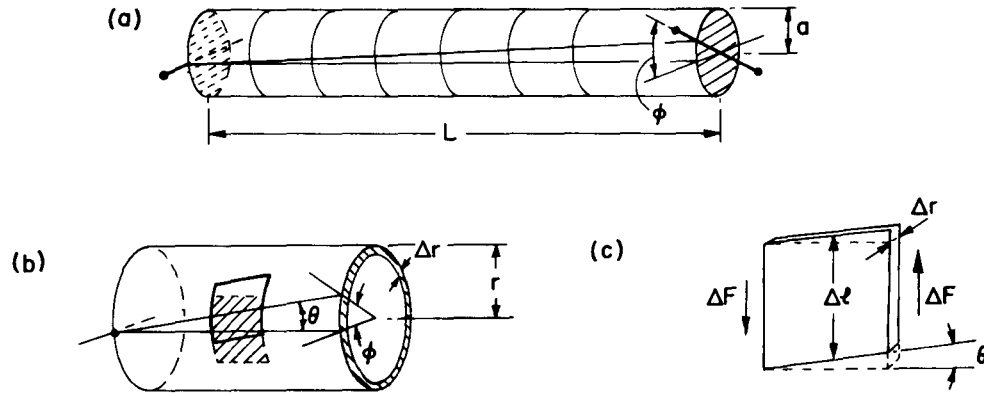


Fig. 38-9. (a) A cylindrical bar in torsion. (b) A cylindrical shell in torsion. (c) Each small piece of the shell is in shear.

Fig. 38-9(a) shows a cylindrical rod of length L , and radius a , with one end twisted by the angle ϕ with respect to the other. If we want to relate the strains to what we already know, we can think of the rod as being made up of many cylindrical shells and work out separately what happens to each shell. We start by looking at a thin, short cylinder of radius r (less than a) and thickness Δr —as drawn in Fig. 38-9(b). Now if we look at a piece of this cylinder that was originally a small square, we see that it has been distorted into a parallelogram. Each such element of the cylinder is in shear, and the shear angle θ is

$$\theta = \frac{r\phi}{L}.$$

The shear stress g in the material is, therefore [from Eq. (38.13)],

$$g = \mu\theta = \mu \frac{r\phi}{L}. \quad (38.21)$$

The shear stress is the tangential force ΔF on the end of the square divided by the area $\Delta l \Delta r$ of the end [see Fig. 38-9(c)]

$$g = \frac{\Delta F}{\Delta l \Delta r}.$$

The force ΔF on the end of such a square contributes a torque $\Delta\tau$ around the axis of the rod equal to

$$\Delta\tau = r \Delta F = rg \Delta l \Delta r. \quad (38.22)$$

The total torque τ is the sum of such torques around a complete circumference of the cylinder. So putting together enough pieces so that the Δl 's add up to $2\pi r$, we find that the total torque, for a *hollow tube*, is

$$rg(2\pi r) \Delta r. \quad (38.23)$$

Or, using (38.21),

$$\tau = 2\pi\mu \frac{r^3 \Delta r \phi}{L}. \quad (38.24)$$

We get that the rotational stiffness, τ/ϕ , of a hollow tube is proportional to the cube of the radius r and to the thickness Δr , and inversely proportional to the length L .

We can now imagine a solid rod to be made up of a series of concentric tubes, each twisted by the same angle ϕ (although the internal *stresses* are different for each tube). The total torque is the sum of the torques required to rotate each shell; for the *solid* rod

$$\tau = 2\pi\mu \frac{\phi}{L} \int r^3 dr,$$

where the integral goes from $r = 0$ to $r = a$, the radius of the rod. Integrating, we have

$$\tau = \mu \frac{\pi a^4}{2L} \phi. \quad (38.25)$$

For a rod in torsion, the torque is proportional to the angle and is proportional to the *fourth power* of the diameter—a rod twice as thick is sixteen times as stiff for torsion.

Before leaving the subject of torsion, let us apply what we have just learned to an interesting problem: torsional waves. If you take a long rod and suddenly twist one end, a wave of twist works its way along the rod, as sketched in Fig. 38-10(a). That's a little more exciting than a steady twist—let's see whether we can work out what happens.

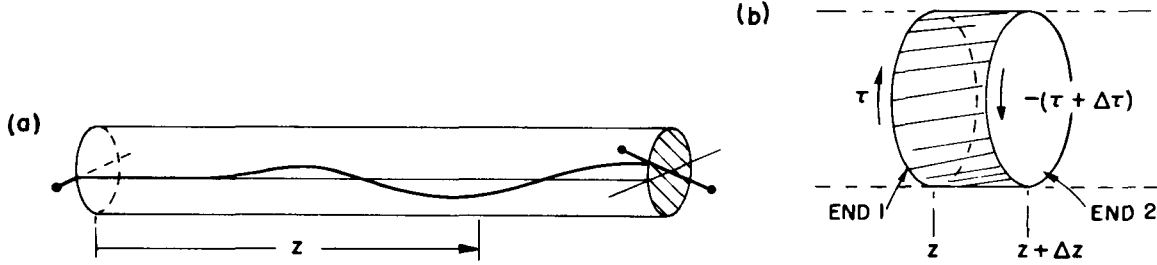


Fig. 38-10. (a) A torsional wave on a rod. (b) A volume element of the rod.

Let z be the distance to some point down the rod. For a static torsion the torque is the same everywhere along the rod, and is proportional to ϕ/L , the total torsion angle over the total length. What matters to the material is the local torsional strain, which is, you will appreciate, $\partial\phi/\partial z$. When the torsion along the rod is not uniform, we should replace Eq. (38.25) by

$$\tau(z) = \mu \frac{\pi a^4}{2} \frac{\partial\phi}{\partial z}. \quad (38.26)$$

Now let's look at what happens to an element of length Δz shown magnified in Fig. 38-10(b). There is a torque $\tau(z)$ at end 1 of the little hunk of rod, and a different torque $\tau(z + \Delta z)$ at end 2. If Δz is small enough, we can use a Taylor expansion and write

$$\tau(z + \Delta z) = \tau(z) + \left(\frac{\partial\tau}{\partial z} \right) \Delta z. \quad (38.27)$$

The net torque $\Delta\tau$ acting on the little piece of rod between z and $z + \Delta z$ is clearly the difference between $\tau(z)$ and $\tau(z + \Delta z)$, or $\Delta\tau = (\partial\tau/\partial z) \Delta z$. Differentiating Eq. (38.26), we get

$$\Delta\tau = \mu \frac{\pi a^4}{2} \frac{\partial^2\phi}{\partial z^2} \Delta z. \quad (38.28)$$

The effect of this net torque is to give an angular acceleration to the little slice of the rod. The mass of the slice is

$$\Delta M = (\pi a^2 \Delta z) \rho,$$

where ρ is the density of the material. We worked out in Chapter 19, Vol. I, that the moment of inertia of a circular cylinder is $mr^2/2$; calling the moment of inertia of our piece ΔI , we have

$$\Delta I = \frac{\pi}{2} \rho a^4 \Delta z. \quad (38.29)$$

Newton's law says the torque is equal to the moment of inertia times the angular acceleration, or

$$\Delta\tau = \Delta I \frac{\partial^2\phi}{\partial t^2}. \quad (38.30)$$

Pulling everything together, we get

$$\mu \frac{\pi a^4}{2} \frac{\partial^2 \phi}{\partial z^2} \Delta z = \frac{\pi}{2} \rho a^4 \Delta z \frac{\partial^2 \phi}{\partial t^2},$$

or

$$\frac{\partial^2 \phi}{\partial z^2} - \frac{\rho}{\mu} \frac{\partial^2 \phi}{\partial t^2} = 0. \quad (38.31)$$

You will recognize this as the one-dimensional wave equation. We have found that waves of torsion will propagate down the rod with the speed

$$C_{\text{shear}} = \sqrt{\frac{\mu}{\rho}}. \quad (38.32)$$

The *denser* the rod—for the same stiffness—the *slower* the waves; and the *stiffer* the rod, the quicker the waves work their way down. The speed does *not* depend upon the diameter of the rod.

Torsional waves are a special example of *shear waves*. In general, shear waves are those in which the strains do not change the *volume* of any part of the material. In torsional waves, we have a particular distribution of such shear stresses—namely, distributed on a circle. But for any arrangement of shear stresses, waves will propagate with the same speed—the one given in Eq. (38.32). For example, the seismologists find such shear waves travelling in the interior of the earth.

We can have another kind of a wave in the elastic world inside a solid material. If you push something, you can start “longitudinal” waves—also called “compressional” waves. They are like the sound waves in air or in water—the displacements are in the same direction as the wave propagation. (At the surfaces of an elastic body there can also be other types of waves—called “Rayleigh waves” or “Love waves.” In them, the strains are neither purely longitudinal nor purely transverse. We will not have time to study them.)

While we’re on the subject of waves, what is the velocity of the pure compressional waves in a *large* solid body like the earth? We say “large” because the speed of sound in a thick body is different from what it is, for instance, along a thin rod. By a “thick” body we mean one in which the transverse dimensions are much larger than the wavelength of the sound. Then, when we push on the object, it cannot expand sideways—it can only compress in one dimension. Fortunately, we have already worked out the special case of the compression of a constrained elastic material. We have also worked out in Chapter 47, Vol. I, the speed of sound waves in a gas. Following the same arguments you can see that the speed of sound in a solid is equal to $\sqrt{Y'/\rho}$, where Y' is the “longitudinal modulus”—or pressure divided by the relative change in length—for the constrained case. This is just the ratio of $\Delta l/l$ to F/A we got in Eq. (38.20). So the speed of the longitudinal waves is given by

$$C_{\text{long}}^2 = \frac{Y'}{\rho} = \frac{1 - \sigma}{(1 + \sigma)(1 - 2\sigma)} \frac{Y}{\rho}. \quad (38.33)$$

So long as σ is between zero and $1/2$, the shear modulus μ is less than Young’s modulus Y , and also Y' is greater than Y , so

$$\mu < Y < Y'.$$

This means that longitudinal waves travel faster than shear waves. One of the most precise ways of measuring the elastic constants of a substance is by measuring the density of the material and the speeds of the two kinds of waves. From this information one can get both Y and σ . It is, incidentally, by measuring the difference in the arrival times of the two kinds of waves from an earthquake that a seismologist can estimate—even from the signals at only one station—the distance to the quake.

38-4 The bent beam

We want now to look at another practical matter—the *bending* of a rod or a beam. What are the forces when we bend a bar of some arbitrary cross section? We will work it out thinking of a bar with a circular cross section, but our answer will be good for any shape. To save time, however, we will cut some corners, so our theory we will work out is only approximate. Our results will be correct only when the radius of the bend is much larger than the thickness of the beam.

Suppose you grab the two ends of a straight bar and bend it into some curve like the one shown in Fig. 38-11. What goes on inside the bar? Well, if it is curved, that means that the material on the inside of the curve is compressed and the material on the outside is stretched. There is some surface which goes along more or less parallel to the axis of the bar that is neither stretched nor compressed. This is called the *neutral surface*. You would expect this surface to be near the “middle” of the cross section. It can be shown (but we won’t do it here) that, for small bending of simple beams, the neutral surface goes through the “center of gravity” of the cross section. This is true only for “pure” bending—if you are not stretching or compressing the beam at the same time.

For pure bending, then, a thin transverse slice of the bar is distorted as shown in Fig. 38-12(a). The material below the neutral surface has a compressional strain which is *proportional to the distance* from the neutral surface; and the material above is stretched, also in proportion to its distance from the neutral surface. So the longitudinal *stretch* Δl is proportional to the height y . The constant of proportionality is just l over the radius of curvature of the bar—see Fig. 38-12:

$$\frac{\Delta l}{l} = \frac{y}{R}.$$

So the force per unit area—the stress—in a small strip at y is also proportional to the distance from the neutral surface

$$\frac{\Delta F}{\Delta A} = Y \frac{y}{R}. \quad (38.34)$$

Now let’s look at the *forces* that would produce such a strain. The forces acting on the little segment drawn in Fig. 38-12 are shown in the figure. If we think of any transverse cut, the forces acting across it are one way above the neutral surface and the other way below. They come in pairs to make a “bending moment” \mathfrak{M} —by which we mean the torque about the neutral line. We can compute the total moment by integrating the force times the distance from the neutral surface for one of the faces of the segment of Fig. 38-12:

$$\mathfrak{M} = \int_{\text{cross sect}} y dF. \quad (38.35)$$

From Eq. (38.34), $dF = Yy/R dA$, so

$$\mathfrak{M} = \frac{Y}{R} \int y^2 dA.$$

The integral of $y^2 dA$ is what we can call the “moment of inertia” of the geometric cross section about a horizontal axis through its “center of mass”;* we will call it I :

$$\mathfrak{M} = \frac{YI}{R} \quad (38.36)$$

$$I = \int y^2 dA. \quad (38.37)$$

* It is, of course, really the moment of inertia of a slice with unit mass per unit area.

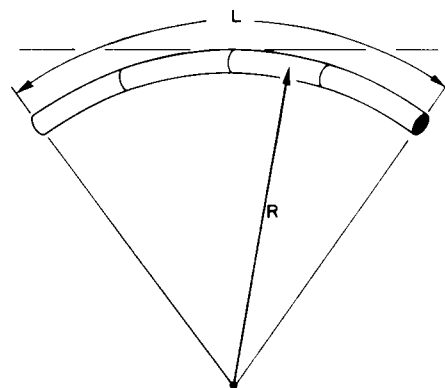


Fig. 38-11. A bent beam.

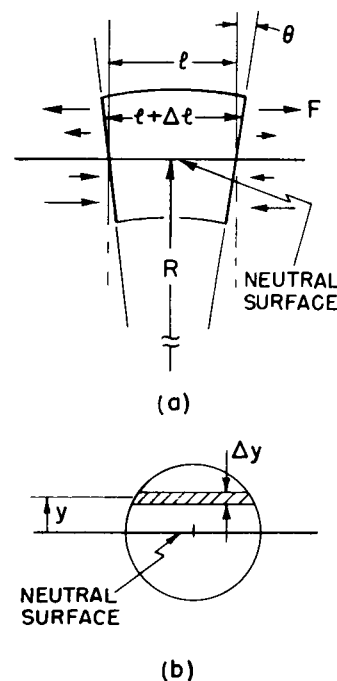


Fig. 38-12. (a) Small segment of a bent beam. (b) Cross section of the beam.

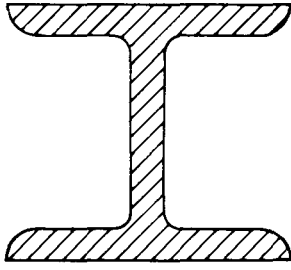


Fig. 38-13. An "I" beam.

Equation (38.36), then, gives us the relation between the bending moment \mathfrak{M} and the curvature $1/R$ of the beam. The "stiffness" of the beam is proportional to Y and to the moment of inertia I . In other words, if you want the stiffest possible beam with a given amount of, say, aluminum, you want to put as much of it as possible as far as you can from the neutral surface, to make a large moment of inertia. You can't carry this to an extreme, however, because then the thing will not curve as we have supposed—it will buckle or twist and become weaker again. But now you see why structural beams are made in the form of an I or an H—as shown in Fig. 38-13.

As an example of the use of our beam equation (38.36), let's work out the deflection of a cantilevered beam with a concentrated force W acting at the free end, as sketched in Fig. 38-14. (By "cantilevered" we simply mean that the beam is supported in such a way that both the position *and* the slope are fixed at one end—it is stuck into a cement wall.) What is the shape of the beam? Let's call the deflection at the distance x from the fixed end z ; we want to know $z(x)$. We'll work it out only for small deflections. We will also assume that the beam is long in comparison with its cross section. Now, as you know from your mathematics courses, the curvature $1/R$ of any curve $z(x)$ is given by

$$\frac{1}{R} = \frac{d^2z/dx^2}{[1 + (dz/dx)^2]^{3/2}}. \quad (38.38)$$

Since we are interested only in small slopes—this is usually the case in engineering structures—we neglect $(dz/dx)^2$ in comparison with 1, and take

$$\frac{1}{R} = \frac{d^2z}{dx^2}. \quad (38.39)$$

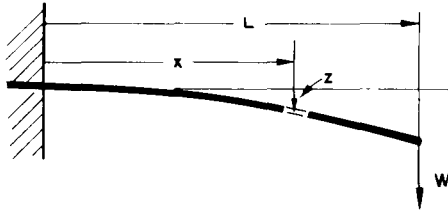


Fig. 38-14. A cantilevered beam with a weight at one end.

We also need to know the bending moment \mathfrak{M} . It is a function of x because it is equal to the torque about the neutral axis of any cross section. Let's neglect the weight of the beam and take only the downward force W at the end of the beam. (You can put in the beam weight yourself if you want.) Then the bending moment at x is

$$\mathfrak{M}(x) = W(L - x),$$

because that is the torque about the point at x , exerted by the weight W —the torque which the beam must support of x . We get

$$W(L - x) = \frac{YI}{R} = YI \frac{d^2z}{dx^2}$$

or

$$\frac{d^2z}{dx^2} = \frac{W}{YI} (L - x). \quad (38.40)$$

This one we can integrate without any tricks; we get

$$z = \frac{W}{YI} \left(\frac{Lx^2}{2} - \frac{x^3}{6} \right), \quad (38.41)$$

using our assumptions that $z(0) = 0$ and that dz/dx is also zero at $x = 0$. That is the shape of the beam. The displacement of the end is

$$z(L) = \frac{W}{YI} \frac{L^3}{3}; \quad (38.42)$$

the displacement of the end of a beam increases as the cube of the length.

In deriving our approximate beam theory, we have assumed that the cross section of the beam did not change when the beam was bent. When the thickness of the beam is small compared to the radius of curvature, the cross section changes very little and our result is O.K. In general, however, this effect cannot be neglected, as you can easily demonstrate for yourselves by bending a soft-rubber eraser in your fingers. If the cross section was originally rectangular, you will find that when

it is bent it bulges at the bottom (see Fig. 38-15). This happens because when we compress the bottom, the material expands sideways—as described by Poisson's ratio. Rubber is easy to bend or stretch, but it is somewhat like a liquid in that it's hard to change the *volume*—as shows up nicely when you bend the eraser. For an incompressible material, Poisson's ratio would be exactly 1/2—for rubber it is nearly that.

38-5 Buckling

We want now to use our beam theory to understand the theory of the “buckling” of beams, or columns, or rods. Consider the situation sketched in Fig. 38-16 in which a rod that would normally be straight is held in its bent shape by two opposite forces that push on the ends of the rod. We would like to calculate the shape of the rod and the *magnitude of the forces* on the ends.

Let the deflection of the rod from the straight line between the ends be $y(x)$, where x is the distance from one end. The bending moment \mathfrak{M} at the point P in the figure is equal to the force F multiplied by the moment arm, which is the perpendicular distance y ,

$$\mathfrak{M}(x) = Fy. \quad (38.43)$$

Using the beam equation (38.36), we have

$$\frac{YI}{R} = Fy. \quad (38.44)$$

For small deflections, we can take $1/R = -d^2y/dx^2$ (the minus sign because the curvature is downward). We get

$$\frac{d^2y}{dx^2} = -\frac{F}{YI} y, \quad (38.45)$$

which is the differential equation of a sine wave. So for *small* deflections, the curve of such a bent beam is a sine curve. The “wavelength” λ of the sine wave is twice the distance L between the ends. If the bending is small, this is just twice the unbent length of the rod. So the curve is

$$y = K \sin \pi x/L.$$

Taking the second derivative, we get

$$\frac{d^2y}{dx^2} = -\frac{\pi^2}{L^2} y.$$

Comparing this to Eq. (38.45), we see that the force is

$$F = \pi^2 \frac{YI}{L^2}. \quad (38.46)$$

For small bendings the force is *independent of the bending displacement y* !

We have, then, the following thing physically. If the force is less than the F given in Eq. (38.46), there will be no bending at all. But if it is slightly *greater* than this force, the material will suddenly bend a large amount—that is, for forces above the critical force $\pi^2 YI/L^2$ (often called the “Euler force”) the beam will “buckle.” If the loading on the second floor of a building exceeds the Euler force for the supporting columns, the building will collapse. Another place where the buckling force is most important is in space rockets. On one hand, the rocket must be able to hold its own weight on the launching pad and endure the stresses during acceleration, on the other hand, it is important to keep the weight of the structure to a minimum, so that the payload and fuel capacity may be made as large as possible.

Actually a beam will not necessarily collapse completely when the force exceeds the Euler force. When the displacements get large, the force is larger than

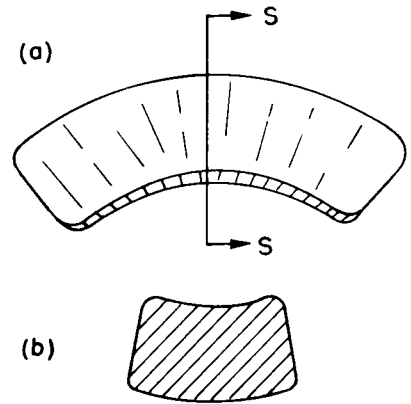


Fig. 38-15 (a) A bent eraser; (b) cross section.

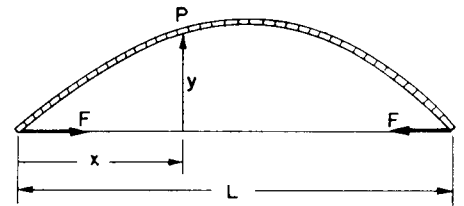


Fig. 38-16. A buckled beam.

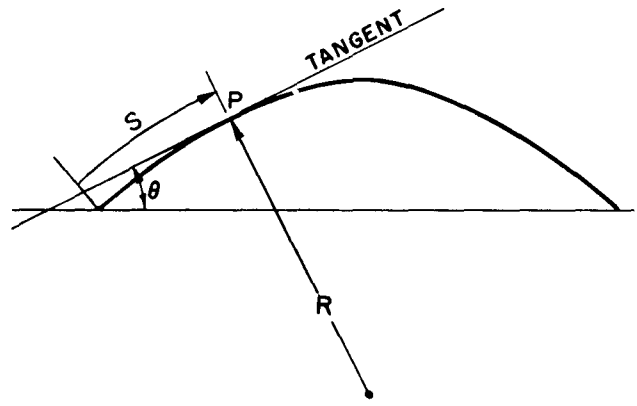


Fig. 38-17. The coordinates S and θ for the curve of a bent beam.

what we have found because of the terms in $1/R$ in Eq. (38.38) that we have neglected. To find the forces for a large bending of the beam, we have to go back to the exact equation, Eq. (38.44), which we had before we used the approximate relation between R and y . Equation (38.44) has a rather simple geometrical property.* It's a little complicated to work out, but rather interesting. Instead of describing the curve in terms of x and y , we can use two new variables: S , the distance along the curve, and θ the slope of the tangent to the curve. See Fig. 38-17. The curvature is the rate of change of angle with distance:

$$\frac{1}{R} = \frac{d\theta}{dS}.$$

We can, therefore write the exact equation (38.44) as

$$\frac{d\theta}{dS} = -\frac{F}{YI} y.$$

If we take the derivative of this equation with respect to S and replace dy/dS by $\sin \theta$, we get

$$\frac{d^2\theta}{dS^2} = -\frac{F}{YI} \sin \theta. \quad (38.47)$$

[If θ is small, we get back Eq. (38.45). Everything is O.K.]

Now it may or may not delight you to know that Eq. (38.47) is exactly the same one you get for the large amplitude oscillations of a pendulum—with F/YI replaced by another constant, of course. We learned way back in Chapter 9, Vol. I, how to find the solution of such an equation by a numerical calculation.† The answers you get are some fascinating curves—known as the curves of the “Elastica.” Figure 38-18 shows three curves for different values of F/YI .

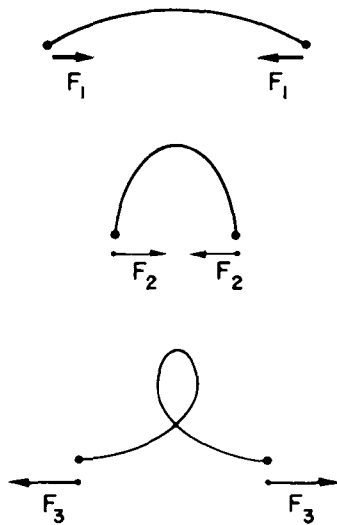


Fig. 38-18. Curves of a bent rod

* The same equation appears, incidentally, in other physical situations—for example, the meniscus at the surface of a liquid contained between parallel planes—and the same geometrical solution can be used.

† The solutions can also be expressed in terms of some functions, called the “Jacobian elliptic functions,” that someone else has already computed.

Elastic Materials

39-1 The tensor of strain

In the last chapter we talked about the distortions of particular elastic objects. In this chapter we want to look at what can happen *in general* inside an elastic material. We would like to be able to describe the conditions of stress and strain inside some big glob of jello which is twisted and squashed in some complicated way. To do this, we need to be able to describe the *local strain* at every point in an elastic body; we can do it by giving a set of six numbers—which are the components of a symmetric tensor—for each point. Earlier, we spoke of the stress tensor (Chapter 31); now we need the tensor of strain.

Imagine that we start with the material initially unstrained and watch the motion of a small speck of “dirt” embedded in the material when the strain is applied. A speck that was at the point P located at $\mathbf{r} = (x, y, z)$ moves to a new position P' at $\mathbf{r}' = (x', y', z')$ as shown in Fig. 39-1. We will call \mathbf{u} the vector displacements from P to P' . Then

$$\mathbf{u} = \mathbf{r}' - \mathbf{r}. \quad (39.1)$$

The displacement \mathbf{u} depends, of course, on which point P we start with, so \mathbf{u} is a vector function of \mathbf{r} —or, if you prefer, of (x, y, z) .

Let's look first at a simple situation in which the strain is constant over the material—so we have what is called a *homogeneous strain*. Suppose, for instance, that we have a block of material and we stretch it uniformly. We just change its dimensions uniformly in one direction—say, in the x -direction, as shown in Fig. 39-2. The motion u_x of a speck at x is proportional to x . In fact,

$$\frac{u_x}{x} = \frac{\Delta l}{l}.$$

We will write u_x this way:

$$u_x = e_{xx}x.$$

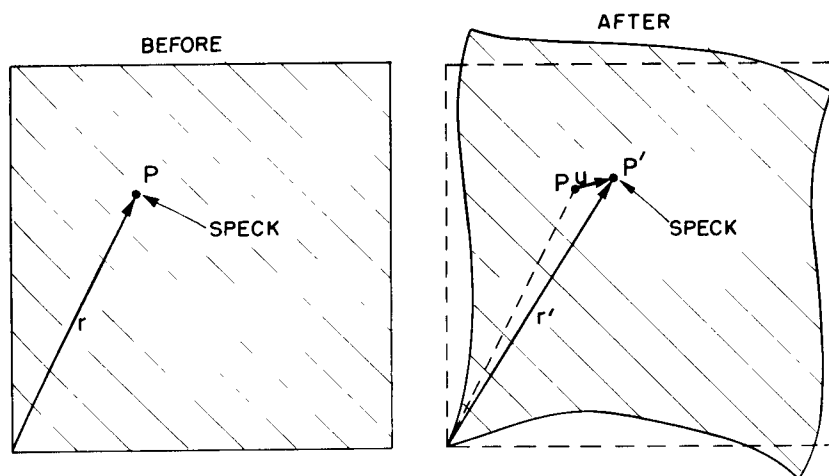


Fig. 39-1. A speck of the material at the point P in an unstrained block moves to P' where the block is strained.

39-1 The tensor of strain

39-2 The tensor of elasticity

39-3 The motions in an elastic body

39-4 Nonelastic behavior

39-5 Calculating the elastic constants

Reference: C. Kittel, *Introduction to Solid State Physics*, John Wiley and Sons, Inc., New York, 2nd ed., 1956.

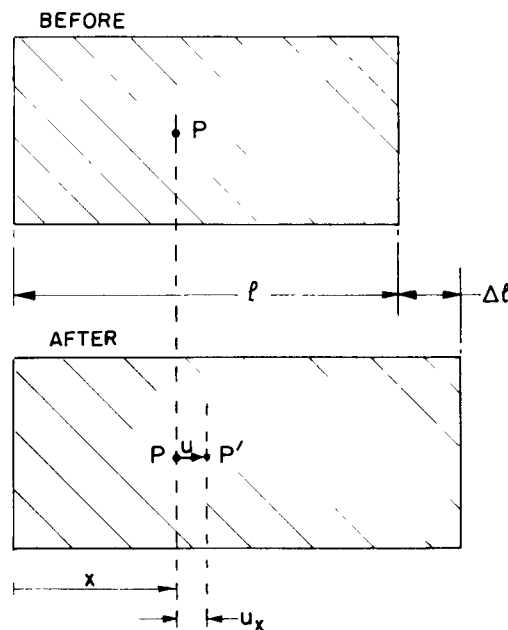


Fig. 39-2. A homogeneous stretch-type strain.

The proportionality constant e_{xx} is, of course, the same thing as $\Delta l/l$. (You will see shortly why we use a double subscript.)

If the strain is not uniform, the relation between u_x and x will vary from place to place in the material. For the general situation, we define the e_{xx} by a kind of local $\Delta l/l$, namely by

$$e_{xx} = \partial u_x / \partial x. \quad (39.2)$$

This number—which is now a function of x , y , and z —describes the amount of stretching in the x -direction throughout the hunk of jello. There may, of course, also be stretching in the y - and z -directions. We describe them by the numbers

$$e_{yy} = \frac{\partial u_y}{\partial y}, \quad e_{zz} = \frac{\partial u_z}{\partial z}. \quad (39.3)$$

We need to be able to describe also the shear-type strains. Suppose we imagine a little cube marked out in the initially undisturbed jello. When the jello is pushed out of shape, this cube may get changed into a parallelogram, as sketched in Fig. 39-3.* In this kind of a strain, the x -motion of each particle is proportional to its y -coordinate,

$$u_x = \frac{\theta}{2} y. \quad (39.4)$$

And there is also a y -motion proportional to x ,

$$u_y = \frac{\theta}{2} x. \quad (39.5)$$

So we can describe such a shear-type strain by writing

$$u_x = e_{xy}y, \quad u_y = e_{yx}x$$

with

$$e_{xy} = e_{yx} = \frac{\theta}{2}.$$

Now you might think that when the strains are not homogeneous we could describe the generalized shear strains by defining the quantities e_{xy} and e_{yx} by

$$e_{xy} = \frac{\partial u_x}{\partial y}, \quad e_{yx} = \frac{\partial u_y}{\partial x}. \quad (39.6)$$

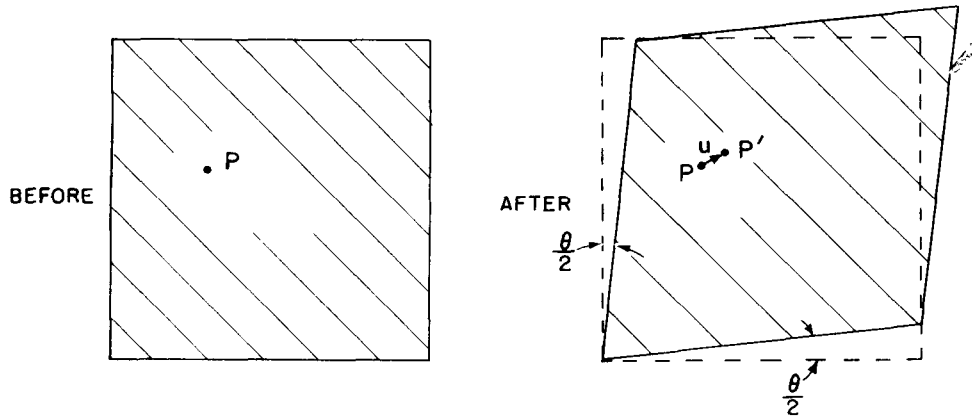


Fig. 39-3. A homogeneous shear strain.

But there is one difficulty. Suppose that the displacements u_x and u_y were given by

$$u_x = \frac{\theta}{2} y, \quad u_y = -\frac{\theta}{2} x.$$

* We choose for the moment to split the total shear angle θ into two equal parts and make the strain symmetric with respect to x and y .

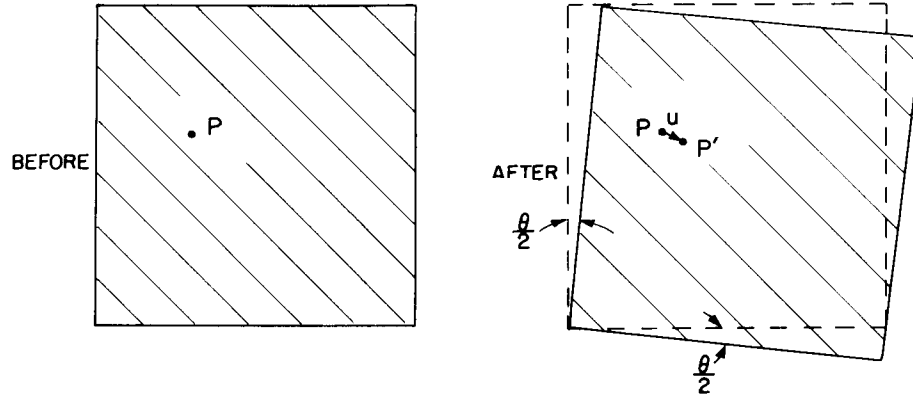


Fig. 39-4. A homogeneous rotation—there is no strain.

They are like Eqs. (39.4) and (39.5) except that the sign of u_y is reversed. With these displacements a little cube in the jello simply gets shifted by the angle $\theta/2$, as shown in Fig. 39-4. There is no strain at all—just a rotation in space. There is no distortion of the material; the *relative* positions of all the atoms are not changed at all. We must somehow make our definitions so that pure rotations are not included in our definitions of a shear strain. The key point is that if $\partial u_y/\partial x$ and $\partial u_x/\partial y$ are equal and opposite, there is no strain; so we can fix things up by *defining*

$$e_{xy} = e_{yx} = \frac{1}{2}(\partial u_y/\partial x + \partial u_x/\partial y).$$

For a pure rotation they are both zero, but for a pure shear we get that e_{xy} is equal to e_{yx} , as we would like.

In the most general distortion—which may include stretching or compression as well as shear—we *define* the state of strain by giving the nine numbers

$$\begin{aligned} e_{xx} &= \frac{\partial u_x}{\partial x}, \\ e_{yy} &= \frac{\partial u_y}{\partial y}, \\ &\vdots \\ e_{xy} &= \frac{1}{2}(\partial u_y/\partial x + \partial u_x/\partial y), \\ &\vdots \end{aligned} \quad (39.7)$$

These are the terms of a *tensor of strain*. Because it is a *symmetric tensor*—our definitions make $e_{xy} = e_{yx}$, always—there are really only six different numbers. You remember (see Chapter 31) that the general characteristic of a tensor is that the terms transform like the products of the components of two vectors. (If \mathbf{A} and \mathbf{B} are vectors, $C_{ij} = A_i B_j$ is a tensor.) Each term of e_{ij} is a product (or the sum of such products) of the components of the vector $\mathbf{u} = (u_x, u_y, u_z)$, and of the operator $\nabla = (\partial/\partial x, \partial/\partial y, \partial/\partial z)$, which we know transforms like a vector. Let's let x_1, x_2 , and x_3 stand for x, y , and z and u_1, u_2 , and u_3 stand for u_x, u_y , and u_z ; then we can write the general term e_{ij} of the strain tensor as

$$e_{ij} = \frac{1}{2}(\partial u_j/\partial x_i + \partial u_i/\partial x_j), \quad (39.8)$$

where i and j can be 1, 2, or 3.

When we have a homogeneous strain—which may include both stretching and shear—all of the e_{ij} are constants, and we can write

$$u_x = e_{xx}x + e_{xy}y + e_{xz}z. \quad (39.9)$$

(We choose our origin of x, y, z at the point where \mathbf{u} is zero.) In this case, the strain tensor e_{ij} gives the relationship between two vectors: the coordinate vector $\mathbf{r} = (x, y, z)$ and the displacement vector $\mathbf{u} = (u_x, u_y, u_z)$.

When the strains are not homogeneous, any piece of the jello may also get somewhat twisted—there will be a local rotation. If the distortions are all small, we would have

$$\Delta u_i = \sum_j (e_{ij} - \omega_{ij}) \Delta x_j, \quad (39.10)$$

where ω_{ij} is an *antisymmetric* tensor,

$$\omega_{ij} = \frac{1}{2}(\partial u_j / \partial x_i - \partial u_i / \partial x_j), \quad (39.11)$$

which describes the rotation. We will, however, not worry any more about rotations, but only about the strains described by the symmetric tensor e_{ij} .

39-2 The tensor of elasticity

Now that we have described the strains, we want to relate them to the internal forces—the stresses in the material. For each small piece of the material, we assume Hooke's law holds and write that the stresses are proportional to the strains. In Chapter 31 we defined the stress tensor S_{ij} as the i th component of the force across a unit area perpendicular to the j -axis. Hooke's law says that each component of S_{ij} is linearly related to *each* of the components of strain. Since S and e each have nine components, there are $9 \times 9 = 81$ possible coefficients which describe the elastic properties of the material. They are constants if the material itself is homogeneous. We write these coefficients as C_{ijkl} and define them by the equation

$$S_{ij} = \sum_{k,l} C_{ijkl} e_{kl}, \quad (39.12)$$

where i, j, k, l all take on the values 1, 2, or 3. Since the coefficients C_{ijkl} relate one tensor to another, they also form a tensor—a tensor of the *fourth rank*. We can call it the *tensor of elasticity*.

Suppose that all the C 's are known and that you put a complicated force on an object of some peculiar shape. There will be all kinds of distortion, and the thing will settle down with some twisted shape. What are the displacements? You can see that it is a complicated problem. If you knew the strains, you could find the stresses from Eq. (39.12)—or vice versa. But the stresses and strains you end up with at any point depend on what happens in all the rest of the material.

The easiest way to get at the problem is by thinking of the energy. When there is a force F proportional to a displacement x , say $F = kx$, the work required for any displacement x is $kx^2/2$. In a similar way, the work w that goes into *each unit volume* of a distorted material turns out to be

$$w = \frac{1}{2} \sum_{ijkl} C_{ijkl} e_{ij} e_{kl}. \quad (39.13)$$

The total work W done in distorting the body is the integral of w over its volume:

$$W = \int \frac{1}{2} \sum_{ijkl} C_{ijkl} e_{ij} e_{kl} d \text{Vol}. \quad (39.14)$$

This is then the potential energy stored in the internal stresses of the material. Now when a body is in equilibrium, this internal energy must be *at a minimum*. So the problem of finding the strains in a body can be solved by finding the set of displacements u throughout the body which will make W a minimum. In Chapter 19 we gave some of the general ideas of the calculus of variations that are used in tackling minimization problems like this. We cannot go into the problem in any more detail here.

What we are mainly interested in now is what we can say about the general properties of the tensor of elasticity. First, it is clear that there are *not* really 81 *different* terms in C_{ijkl} . Since both S_{ij} and e_{ij} are symmetric tensors, each with only six different terms, there can be at most 36 different terms in C_{ijkl} . There are, however, usually many fewer than this.

Let's look at the special case of a cubic crystal. In it, the energy density w starts out like this:

$$w = \frac{1}{2} \{ C_{xxxx} e_{xx}^2 + C_{xxxy} e_{xx} e_{xy} + C_{xxxz} e_{xx} e_{xz} + C_{xxyx} e_{xx} e_{xy} + C_{xyyy} e_{xx} e_{yy} + \dots \text{etc} \dots + C_{yyyy} e_{yy}^2 + \dots \text{etc} \dots \text{etc} \dots \}, \quad (39.15)$$

with 81 terms in all! Now a cubic crystal has certain symmetries. In particular, if the crystal is rotated 90° , it has the same physical properties. It has the same stiffness for stretching in the y -direction as for stretching in the x -direction. Therefore, if we change our definition of the coordinate directions x and y in Eq. (39.15), the energy wouldn't change. It must be that for a cubic crystal

$$C_{xxxx} = C_{yyyy} = C_{zzzz}. \quad (39.16)$$

Next we can show that the terms like C_{xxxy} must be zero. A cubic crystal has the property that it is symmetric under a *reflection* about any plane perpendicular to one of the axes. If we replace y by $-y$, nothing is different. But changing y to $-y$ changes e_{xy} to $-e_{xy}$ —a displacement which was toward $+y$ is now toward $-y$. If the energy is not to change, C_{xxxy} must go into $-C_{xxxy}$ when we make a reflection. But a reflected crystal is the same as before, so C_{xxxy} must be the *same* as $-C_{xxxy}$. This can happen only if both are zero.

You say, "But the same argument will make $C_{yyyy} = 0$!" No, because there are *four* y 's. The sign changes once for each y , and four minuses make a plus. If there are *two* or *four* y 's, the term does not have to be zero. It is zero only when there is *one*, or *three*. So, for a cubic crystal, any nonzero term of C will have only an *even number* of identical subscripts. (The arguments we have made for y obviously hold also for x and z .) We might then have terms like C_{xxyy} , C_{xyxy} , C_{yyxx} , and so on. We have already shown, however, that if we change all x 's to y 's and *vice versa* (or all z 's and x 's, and so on) we must get—for a cubic crystal—the same number. This means that there are *only three different* nonzero possibilities:

$$\begin{aligned} C_{xxxx} (= C_{yyyy} = C_{zzzz}), \\ C_{xxyy} (= C_{yyxx} = C_{zzzz}, \text{ etc.}), \\ C_{xyxy} (= C_{yxxy} = C_{zzzz}, \text{ etc.}). \end{aligned} \quad (39.17)$$

For a cubic crystal, then, the energy density will look like this:

$$w = \frac{1}{2} \{ C_{xxxx} (e_{xx}^2 + e_{yy}^2 + e_{zz}^2) + 2C_{xxyy} (e_{xx} e_{yy} + e_{yy} e_{zz} + e_{zz} e_{xx}) + 4C_{xyxy} (e_{xy}^2 + e_{yz}^2 + e_{zx}^2) \}. \quad (39.18)$$

For an isotropic—that is, noncrystalline—material, the symmetry is still higher. The C 's must be the same for *any* choice of the coordinate system. Then it turns out that there is another relation among the C 's, namely, that

$$C_{xxxx} = C_{xxyy} + C_{xyxy}. \quad (39.19)$$

We can see that this is so by the following general argument. The stress tensor S_{ij} has to be related to e_{ij} in a way that doesn't depend at all on the coordinate directions—it must be related only by *scalar* quantities. "That's easy," you say. "The only way to obtain S_{ij} from e_{ij} is by multiplication by a scalar constant. It's just Hooke's law. It must be that $S_{ij} = (\text{const})e_{ij}$." But that's not quite right; there could also be the *unit tensor* δ_{ij} multiplied by some scalar, linearly related to e_{ij} . The only invariant you can make that is linear in the e 's is $\sum e_{ii}$ (It transforms like $x^2 + y^2 + z^2$, which is a scalar.) So the most general form for the equation relating S_{ij} to e_{ij} —for isotropic materials—is

$$S_{ij} = 2\mu e_{ij} + \lambda \left(\sum_k e_{kk} \right) \delta_{ij}. \quad (39.20)$$

(The first constant is usually written as *two* times μ ; then the coefficient μ is equal

to the shear modulus we defined in the last chapter.) The constants μ and λ are called the Lamé elastic constants. Comparing Eq. (39.20) with Eq. (39.12), you see that

$$\begin{aligned} C_{xxyy} &= \lambda, \\ C_{xyxy} &= 2\mu, \\ C_{xxxx} &= 2\mu + \lambda. \end{aligned} \quad (39.21)$$

So we have proved that Eq. (39.19) is indeed true. You also see that the elastic properties of an isotropic material are completely given by two constants, as we said in the last chapter.

The C 's can be put in terms of any two of the elastic constants we have used earlier—for instance, in terms of Young's modulus Y and Poisson's ratio σ . We will leave it for you to show that

$$\begin{aligned} C_{xxxx} &= \frac{Y}{1+\sigma} \left(1 + \frac{\sigma}{1-2\sigma} \right), \\ C_{xxyy} &= \frac{Y}{1+\sigma} \left(\frac{\sigma}{1-2\sigma} \right), \\ C_{xyxy} &= \frac{Y}{(1+\sigma)}. \end{aligned} \quad (39.22)$$

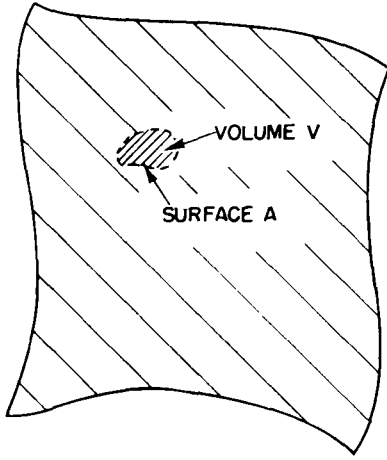


Fig. 39-5. A small volume element V bounded by the surface A .

39-3 The motions in an elastic body

We have pointed out that for an elastic body *in equilibrium* the internal stresses adjust themselves to make the energy a minimum. Now we take a look at what happens when the internal forces are *not* in equilibrium. Let's say we have a small piece of the material inside some surface A . See Fig. 39-5. If the piece is in equilibrium, the total force F acting on it must be zero. We can think of this force as being made up of two parts. There could be one part due to "external" forces like gravity, which act from a distance on the matter in the piece to produce a *force per unit volume* f_{ext} . The total external force F_{ext} is the integral of f_{ext} over the volume of the piece:

$$F_{\text{ext}} = \int f_{\text{ext}} dV. \quad (39.23)$$

In equilibrium, this force would be balanced by the total force F_{int} from the neighboring material which acts across the surface A . When the piece is *not* in equilibrium—if it is moving—the sum of the internal and external forces is equal to the mass times the acceleration. We would have

$$F_{\text{ext}} + F_{\text{int}} = \int \rho \mathbf{r} dV, \quad (39.24)$$

where ρ is the density of the material, and \mathbf{r} is its acceleration. We can now combine Eqs. (39.23) and (39.24), writing

$$F_{\text{int}} = \int_v (-f_{\text{ext}} + \rho \mathbf{r}) dV. \quad (39.25)$$

We will simplify our writing by defining

$$\mathbf{f} = -f_{\text{ext}} + \rho \mathbf{r}. \quad (39.26)$$

Then Eq. (39.25) is written

$$F_{\text{int}} = \int_v \mathbf{f} dV. \quad (39.27)$$

What we have called F_{int} is related to the stresses in the material. The stress tensor S_{ij} was defined (Chapter 31) so that the x -component of the force dF across a surface element da , whose unit normal is \mathbf{n} , is given by

$$dF_x = (S_{xx}n_x + S_{xy}n_y + S_{xz}n_z) da. \quad (39.28)$$

The x -component of \mathbf{F}_{int} on our little piece is then the integral of dF_x over the surface. Substituting this into the x -component of Eq. (39.27), we get

$$\int_A (S_{xx}n_x + S_{xy}n_y + S_{xz}n_z) da = \int_v f_x dV. \quad (39.29)$$

We have a surface integral related to a volume integral—and that reminds us of something we learned in electricity. Note that if you ignore the first subscript x on each of the S 's in the left-hand side of Eq. (39.29), it looks just like the integral of a quantity " \mathbf{S} " $\cdot \mathbf{n}$ —that is, the normal component of a vector—over the surface. It would be the flux of " \mathbf{S} " out of the volume. And this could be written, using Gauss law, as the volume integral of the divergence of " \mathbf{S} ". It is, in fact, true whether the x -subscript is there or not—it is just a mathematical theorem you get by integrating by parts. In other words, we can change Eq. (39.29) into

$$\int_v \left(\frac{\partial S_{xx}}{\partial x} + \frac{\partial S_{xy}}{\partial y} + \frac{\partial S_{xz}}{\partial z} \right) dV = \int_v f_x dV \quad (39.30)$$

Now we can leave off the volume integrals and write the differential equation for the general component of \mathbf{f} as

$$f_i = \sum_j \frac{\partial S_{ij}}{\partial x_j}. \quad (39.31)$$

This tells us how the force per unit volume is related to the stress tensor S_{ij} .

The theory of the motions inside a solid works this way. If we start out knowing the initial displacements—given by, say, \mathbf{u} —we can work out the strains e_{ij} . From the strains we can get the stresses from Eq. (39.12). From the stresses we can get the force density \mathbf{f} in Eq. (39.31). Knowing \mathbf{f} , we can get, from Eq. (39.26), the acceleration \mathbf{r} of the material, which tells us how the displacements will be changing. Putting everything together, we get the horrible equation of motion for an elastic solid. We will just write down the results that come out for an isotropic material. If you use (39.20) for S_{ij} , and write the e_{ij} as $\frac{1}{2}\partial u_i/\partial x_j + \partial u_j/\partial x_i$, you end up with the vector equation

$$\mathbf{f} = (\lambda + \mu) \nabla(\nabla \cdot \mathbf{u}) + \mu \nabla^2 \mathbf{u}. \quad (39.32)$$

You can, in fact, see that the equation relating \mathbf{f} and \mathbf{u} *must* have this form. The force must depend on the second derivatives of the displacements \mathbf{u} . What second derivatives of \mathbf{u} are there that are vectors? One is $\nabla(\nabla \cdot \mathbf{u})$; that's a true vector. The only other one is $\nabla^2 \mathbf{u}$. So the most general form is

$$\mathbf{f} = a \nabla(\nabla \cdot \mathbf{u}) + b \nabla^2 \mathbf{u},$$

which is just (39.32) with a different definition of the constants. You may be wondering why we don't have a third term using $\nabla \times \nabla \times \mathbf{u}$, which is also a vector. But remember that $\nabla \times \nabla \times \mathbf{u}$ is the same thing as $\nabla^2 \mathbf{u} - \nabla(\nabla \cdot \mathbf{u})$, so it is a linear combination of the two terms we have. Adding it would add nothing new. We have proved once more that isotropic material has only two elastic constants.

For the equation of motion of the material, we can set (39.32) equal to $\rho \partial^2 \mathbf{u}/\partial t^2$ —neglecting for now any body forces like gravity—and get

$$\rho \frac{\partial^2 \mathbf{u}}{\partial t^2} = (\lambda + \mu) \nabla(\nabla \cdot \mathbf{u}) + \mu \nabla^2 \mathbf{u}. \quad (39.33)$$

It looks something like the wave equation we had in electromagnetism, except that there is an additional complicating term. For materials whose elastic properties are everywhere the same we can see what the general solutions look like in the following way. You will remember that any vector field can be written as the sum of two vectors: one whose divergence is zero, and the other whose curl is zero. In

other words, we can put

$$\mathbf{u} = \mathbf{u}_1 + \mathbf{u}_2, \quad (39.34)$$

where

$$\nabla \cdot \mathbf{u}_1 = 0, \quad \nabla \times \mathbf{u}_2 = 0. \quad (39.35)$$

Substituting $\mathbf{u}_1 + \mathbf{u}_2$ for \mathbf{u} in (39.33), we get

$$\rho \partial^2 / \partial t^2 [\mathbf{u}_1 + \mathbf{u}_2] = (\lambda + \mu) \nabla (\nabla \cdot \mathbf{u}_2) + \mu \nabla^2 (\mathbf{u}_1 + \mathbf{u}_2). \quad (39.36)$$

We can eliminate \mathbf{u}_1 by taking the divergence of this equation,

$$\rho \partial^2 / \partial t^2 (\nabla \cdot \mathbf{u}_2) = (\lambda + \mu) \nabla^2 (\nabla \cdot \mathbf{u}_2) + \mu \nabla \cdot \nabla^2 \mathbf{u}_2.$$

Since the operators (∇^2) and $(\nabla \cdot)$ can be interchanged, we can factor out the divergence to get

$$\nabla \cdot \{ \rho \partial^2 \mathbf{u}_2 / \partial t^2 - (\lambda + 2\mu) \nabla^2 \mathbf{u}_2 \} = 0. \quad (39.37)$$

Since $\nabla \times \mathbf{u}_2$ is zero by definition, the curl of the bracket $\{ \}$ is also zero; so the bracket itself is identically zero, and

$$\rho \partial^2 \mathbf{u}_2 / \partial t^2 = (\lambda + 2\mu) \nabla^2 \mathbf{u}_2. \quad (39.38)$$

This is the vector wave equation for waves which move at the speed $C_2 = \sqrt{(\lambda + 2\mu)/\rho}$. Since the curl of \mathbf{u}_2 is zero, there is no shearing associated with this wave; this wave is just the compressional—sound-type—wave we discussed in the last chapter, and the velocity is just what we found for C_{long} .

In a similar way—by taking the curl of Eq. (39.36)—we can show that \mathbf{u}_1 satisfies the equation

$$\rho \partial^2 \mathbf{u}_1 / \partial t^2 = \mu \nabla^2 \mathbf{u}_1. \quad (39.39)$$

This is again a vector wave equation for waves with the speed $C_2 = \sqrt{\mu/\rho}$. Since $\nabla \cdot \mathbf{u}_1$ is zero, \mathbf{u}_1 produces no changes in density; the vector \mathbf{u}_1 corresponds to the transverse, or shear-type, wave we saw in the last chapter, and $C_2 = C_{\text{shear}}$.

If we wished to know the static stresses in an isotropic material, we could, in principle, find them by solving Eq. (39.32) with f equal to zero—or equal to the static body forces from gravity such as ρg —under certain conditions which are related to the forces acting on the surfaces of our large block of material. This is somewhat more difficult to do than the corresponding problems in electromagnetism. It is more difficult, first, because the equations are a little more difficult to handle, and second, because the shape of the elastic bodies we are likely to be interested in are usually much more complicated. In electromagnetism, we are often interested in solving Maxwell's equations around relatively simple geometric shapes such as cylinders, spheres, and so on, since these are convenient shapes for electrical devices. In elasticity, the objects we would like to analyze may have quite complicated shapes—like a crane hook, or an automobile crankshaft, or the rotor of a gas turbine. Such problems can sometimes be worked out approximately by numerical methods, using the minimum energy principle we mentioned earlier. Another way is to use a model of the object and measure the internal strains experimentally, using polarized light.

It works this way: When a transparent isotropic material—for example, a clear plastic like lucite—is put under stress, it becomes birefringent. If you put polarized light through it, the plane of polarization will be rotated by an amount related to the stress: by measuring the rotation, you can measure the stress. Figure 39-6 shows how such a setup might look. Figure 39-7 is a photograph of a photoelastic model of a complicated shape under stress.

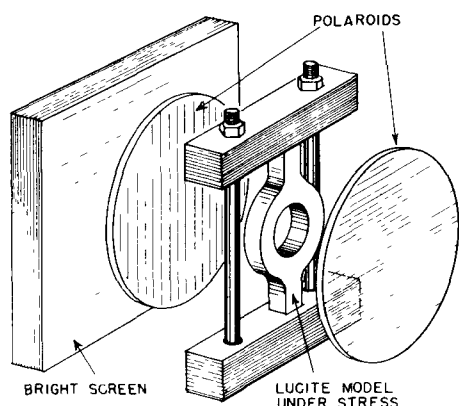


Fig. 39-6. Measuring internal stresses with polarized light.

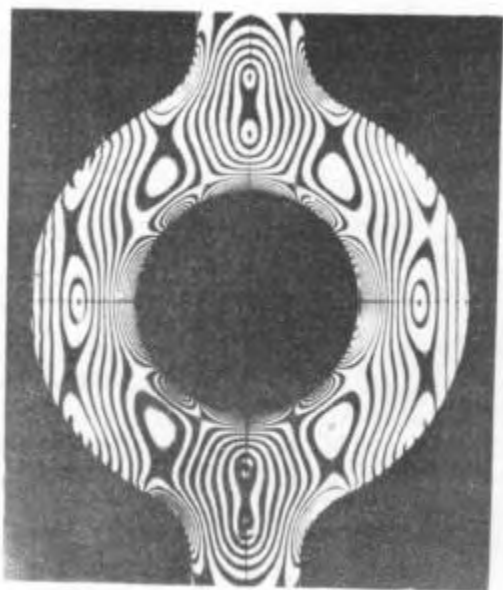


Fig. 39-7. A stressed plastic model as seen between crossed polaroids. [From F. W. Sears, *Optics*, Addison-Wesley Publishing Co., Reading, Mass., 1949.]

39-4 Nonelastic behavior

In all that has been said so far, we have assumed that stress is proportional to strain; in general, that is *not* true. Figure 39-8 shows a typical stress-strain curve for a ductile material. For small strains, the stress is proportional to the

strain. Eventually, however, after a certain point, the relationship between stress and strain begins to deviate from a straight line. For many materials—the ones we would call “brittle”—the object breaks for strains only a little above the point where the curve starts to bend over. In general, there are other complications in the stress-strain relationship. For example, if you strain an object, the stresses may be high at first, but decrease slowly with time. Also if you go to high stresses, but still not to the “breaking” point, when you lower the strain the stress will return along a different curve. There is a small hysteresis effect (like the one we saw between B and H in magnetic materials).

The stress at which a material will break varies widely from one material to another. Some materials will break when the maximum *tensile* stress reaches a certain value. Other materials will fail when the maximum *shear* stress reaches a certain value. Chalk is an example of a material which is much weaker in tension than in shear. If you pull on the ends of a piece of blackboard chalk, the chalk will break perpendicular to the direction of the applied stress, as shown in Fig. 39-9(a). It breaks perpendicular to the applied force because it is only a bunch of particles packed together which are easily pulled apart. The material is, however, much harder to shear, because the particles get in each other's way. Now you will remember that when we had a rod in torsion there was a shear all around it. Also, we showed that a shear was equivalent to a combination of a tension and compression at 45° . For these reasons, if you *twist* a piece of blackboard chalk, it will break along a complicated surface which starts out at 45° to the axis. A photograph of a piece of chalk broken in this way is shown in Fig. 39-9(b). The chalk breaks where the material is in maximum tension.

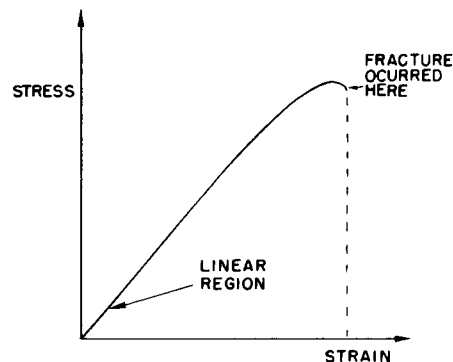


Fig. 39-8. A typical stress-strain relation for large strains.

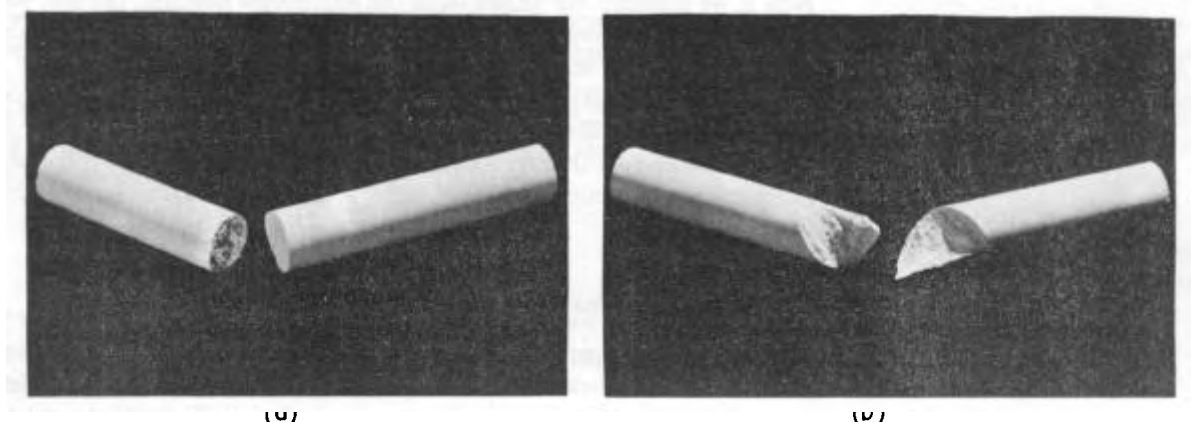


Fig. 39-9. (a) A piece of chalk broken by pulling on the ends; (b) a piece broken by twisting.

Other materials behave in strange and complicated ways. The more complicated the materials are, the more interesting their behavior. If we take a sheet of “Saran-Wrap” and crumple it up into a ball and throw it on the table, it slowly unfolds itself and returns toward its original flat form. At first sight, we might be tempted to think that it is inertia which prevents it from returning to its original form. However, a simple calculation shows that the inertia is several orders of magnitude too small to account for the effect. There appear to be two important competing effects: “something” inside the material “remembers” the shape it had initially and “tries” to get back there, but something else “prefers” the new shape and “resists” the return to the old shape.

We will not attempt to describe the mechanism at play in the Saran plastic, but you can get an idea of how such an effect might come about from the following *model*. Suppose you imagine a material made of long, flexible, but strong, fibers mixed together with some hollow cells filled with a viscous liquid. Imagine also that there are narrow pathways from one cell to the next so the liquid can leak slowly from a cell to its neighbor. When we crumple a sheet of this stuff, we distort the long fibers, squeezing the liquid out of the cells in one place and forcing it into other cells which are being stretched. When we let go, the long fibers try to

return to their original shape. But to do this, they have to force the liquid back to its original location—which will happen relatively slowly because of the viscosity. The forces we apply in crumpling the sheet are much larger than the forces exerted by the fibers. We can crumple the sheet quickly, but it will return more slowly. It is undoubtedly a combination of large stiff molecules and smaller, movable ones in the Saran-Wrap that is responsible for its behavior. This idea also fits with the fact that the material returns more quickly to its original shape when it's warmed up than when it's cold—the heat increases the mobility (decreases the viscosity) of the smaller molecules.

Although we have been discussing how Hooke's law breaks down, the remarkable thing is perhaps not that Hooke's law breaks down for large strains but that it should be so generally true. We can get some idea of why this might be by looking at the strain energy in a material. To say that the stress is proportional to the strain is the same thing as saying that the strain energy varies as the square of the strain. Suppose we have a rod and we twist it through a small angle θ . If Hooke's law holds, the strain energy should be proportional to the square of θ . Suppose we were to assume that the energy were some arbitrary function of the angle; we could write it as a Taylor expansion about zero angle

$$U(\theta) = U(0) + U'(0)\theta + \frac{1}{2}U''(0)\theta^2 + \frac{1}{6}U'''(0)\theta^3 \dots \quad (39.40)$$

The torque τ is the derivative of U with respect to angle; we would have

$$\tau(\theta) = U'(0) + U''(0)\theta + \frac{1}{2}U'''(0)\theta^2 + \dots \quad (39.41)$$

Now if we measure our angles from the *equilibrium* position, the first term is zero. So the first remaining term is proportional to θ ; and for small enough angles, it will dominate the term in θ^2 . [Actually, materials are sufficiently symmetric internally so that $\tau(\theta) = -\tau(-\theta)$; the term in θ^2 will be zero, and the departures from linearity would come only from the θ^3 term. There is, however, no reason why this should be true for compressions and tensions.] The thing we have not explained is why materials usually break soon after the higher-order terms become significant.

39-5 Calculating the elastic constants

As our last topic on elasticity we would like to show how one could try to calculate the elastic constants of a material, starting with some knowledge of the properties of the atoms which make up the material. We will take only the simple case of an *ionic* cubic crystal like sodium chloride. When a crystal is strained, its volume or its shape is changed. Such changes result in an increase in the potential energy of the crystal. To calculate the change in strain energy, we have to know where each atom goes. In complicated crystals, the atoms will rearrange themselves in the lattice in very complicated ways to make the total energy as small as possible. This makes the computation of the strain energy rather difficult. In the case of a simple cubic crystal, however, it is easy to see what will happen. The distortions inside the crystal will be geometrically similar to the distortions of the outside boundaries of the crystal.

We can calculate the elastic constants for a cubic crystal in the following way. First, we assume some force law between each pair of atoms in the crystal. Then, we calculate the change in the internal energy of the crystal when it is distorted from its equilibrium shape. This gives us a relation between the energy and the strains which is quadratic in all the strains. Comparing the energy obtained this way with Eq. (39.13), we can identify the coefficient of each term with the elastic constants C_{ijkl} .

For our example we will assume a simple force law: that the force between neighboring atoms is a *central* force, by which we mean that it acts along the line between the two atoms. We would expect the forces in ionic crystals to be like this, since they are just primarily Coulomb forces. (The forces of covalent bonds are usually more complicated, since they can exert a sideways push on a nearby

atom; we will leave out this complication.) We are also going to include only the forces between each atom and its *nearest* and *next-nearest* neighbors. In other words, we will make an approximation which neglects all forces beyond the next-nearest neighbor. The forces we will include are shown for the xy -plane in Fig. 39-10(a). The corresponding forces in the yz - and zx -planes also have to be included.

Since we are only interested in the elastic coefficients which apply to small strains, and therefore only want the terms in the energy which vary quadratically with the strains, we can imagine that the force between each atom pair varies linearly with the displacements. We can then imagine that each pair of atoms is joined by a linear spring, as drawn in Fig. 39-10(b). All of the springs between a sodium atom and a chlorine atom should have the same spring constant, say k_1 . The springs between two sodiums and between two chlorines could have different constants, but we will make our discussion simpler by taking them equal; we call them k_2 . (We could come back later and make them different after we have seen how the calculations go.)

Now we assume that the crystal is distorted by a homogeneous strain described by the strain tensor e_{ij} . In general, it will have components involving x , y , and z ; but we will consider now only a strain with the three components e_{xx} , e_{xy} , and e_{yy} so that it will be easy to visualize. If we pick one atom as our origin, the displacement of every other atom is given by equations like Eq. (39.9):

$$\begin{aligned} u_x &= e_{xx}x + e_{xy}y, \\ u_y &= e_{xy}x + e_{yy}y. \end{aligned} \quad (39.42)$$

Suppose we call the atom at $x = y = 0$ "atom 1" and number its neighbors in the xy -plane as shown in Fig. 39-11. Calling the lattice constant a , we get the x and y displacements u_x and u_y listed in Table 39-1.

Now we can calculate the energy stored in the springs, which is $k^2/2$ times the square of the extension for each spring. For example, the energy in the horizontal spring between atom 1 and atom 2 is

$$\frac{k_1(e_{xx}a)^2}{2}. \quad (39.43)$$

Note that to first order, the y -displacement of atom 2 does not change the length of the spring between atom 1 and atom 2. To get the strain energy in a diagonal spring, such as that to atom 3, however, we need to calculate the change in length due to both the horizontal and vertical displacements. For small displacements from the

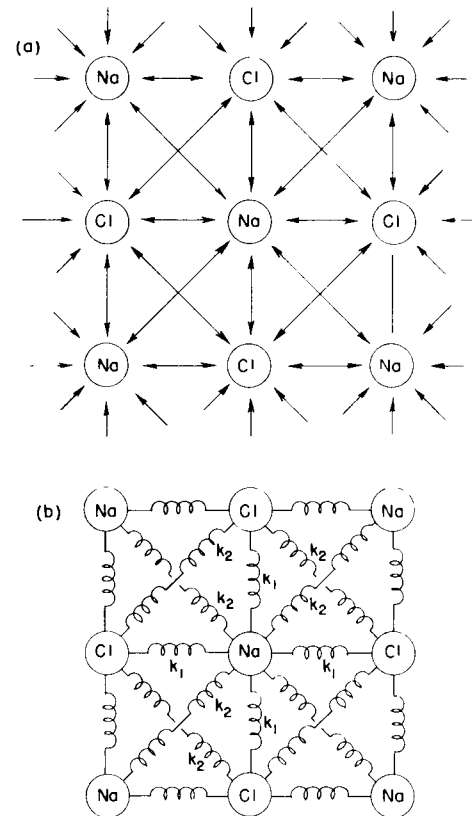


Fig. 39-10. (a) The interatomic forces we are taking into account; (b) a model in which the atoms are connected by springs.

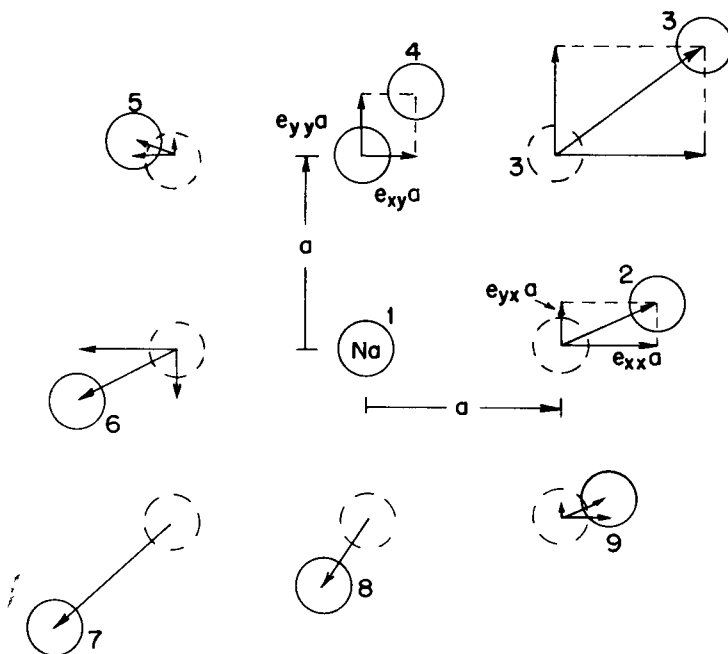


Fig. 39-11. The displacements of the nearest and next-nearest neighbors of atom 1 (exaggerated).

Table 39-1

Atom	Location x, y	u_x	u_y	k
1	0, a	0	0	—
2	$a, 0$	$e_{xx}a$	$e_{yx}a$	k_1
3	a, a	$(e_{xx} + e_{xy})a$	$(e_{yx} + e_{yy})a$	k_2
4	0, a	$e_{xy}a$	$e_{yy}a$	k_1
5	$-a, a$	$(-e_{xx} + e_{xy})a$	$(-e_{yx} + e_{yy})a$	k_2
6	$-a, 0$	$-e_{xx}a$	$-e_{yx}a$	k_1
7	$-a, -a$	$-(e_{xx} + e_{xy})a$	$-(e_{yx} + e_{yy})a$	k_2
8	0, $-a$	$-e_{xy}a$	$-e_{yy}a$	k_1
9	$a, -a$	$(e_{xx} - e_{xy})a$	$(e_{yx} - e_{yy})a$	k_2

original cube, we can write the change in the distance to atom 3 as the sum of the components of u_x and u_y in the diagonal direction, namely as

$$\frac{1}{\sqrt{2}} (u_x + u_y).$$

Using the values of u_x and u_y from the table, we get the energy

$$\frac{k_2}{2} \left(\frac{u_x + u_y}{\sqrt{2}} \right)^2 = \frac{k_2 a^2}{4} (e_{xx} + e_{yx} + e_{xy} + e_{yy})^2. \quad (39.44)$$

For the total energy for all the springs in the xy -plane, we need the sum of eight terms like (39.43) and (39.44). Calling this energy U_0 , we get

$$\begin{aligned} U_0 = \frac{a^2}{2} \left\{ k_1 e_{xx}^2 + \frac{k_2}{2} (e_{xx} + e_{yx} + e_{xy} + e_{yy})^2 \right. \\ + k_1 e_{yy}^2 + \frac{k_2}{2} (e_{xx} - e_{yx} - e_{xy} + e_{yy})^2 \\ + k_1 e_{xx}^2 + \frac{k_2}{2} (e_{xx} + e_{yx} + e_{xy} + e_{yy})^2 \\ \left. + k_1 e_{yy}^2 + \frac{k_2}{2} (e_{xx} - e_{yx} - e_{xy} + e_{yy})^2 \right\}. \quad (39.45) \end{aligned}$$

To get the total energy of all the springs connected to atom 1, we must make one addition to the energy in Eq. (39.45). Even though we have only x - and y -components of the strain, there are still some energies associated with the next-nearest neighbors off the xy -plane. This additional energy is

$$k_2 (e_{xx}^2 a^2 + e_{yy}^2 a^2). \quad (39.46)$$

The elastic constants are related to the energy density w by Eq. (39.13). The energy we have calculated is the energy associated with one atom, or rather, it is *twice* the energy per atom, since one-half of the energy of each spring should be assigned to each of the two atoms it joins. Since there are $1/a^3$ atoms per unit volume, w and U_0 are related by

$$w = \frac{U_0}{2a^3}.$$

To find the elastic constants C_{ijkl} , we need only to expand out the squares in Eq. (39.45)—adding the terms of (39.46)—and compare the coefficients of $e_{ij}e_{kl}$ with the corresponding coefficient in Eq. (39.13). For example, collecting the terms

in e_{xx}^2 and in e_{yy}^2 , we get the factor

$$(k_1 + 2k_2)a^2,$$

so

$$C_{xxxx} = C_{yyyy} = \frac{k_1 + 2k_2}{a}.$$

For the remaining terms, there is a slight complication. Since we cannot distinguish the product of two terms like $e_{xx}e_{yy}$ from $e_{yy}e_{xx}$, the coefficient of such terms in our energy is equal to the sum of two terms in Eq. (39.13). The coefficient of $e_{xx}e_{yy}$ in Eq. (39.45) is $2k_2$, so we have that

$$(C_{xxyy} + C_{yyxx}) = \frac{2k_2}{a}.$$

But because of the symmetry in our crystal, $C_{xxyy} = C_{yyxx}$, so we have that

$$C_{xxyy} = C_{yyxx} = \frac{k_2}{a}.$$

By a similar process, we can also get

$$C_{xyxy} = C_{yxxy} = \frac{k_2}{a}.$$

Finally, you will notice that any term which involves either x or y only once is zero—as we concluded earlier from symmetry arguments. Summarizing our results:

$$\begin{aligned} C_{xxxx} &= C_{yyyy} = \frac{k_1 + 2k_2}{a}, \\ C_{xxyy} &= C_{yyxx} = \frac{k_2}{a}, \\ C_{xyxy} &= C_{yxxy} = \frac{k_2}{a}, \\ C_{xxxy} &= C_{xyyy} = \text{etc.} = 0. \end{aligned} \quad (39.47)$$

We have been able to relate the bulk elastic constants to the atomic properties which appear in the constants k_1 and k_2 . In our particular case, $C_{xxyy} = C_{yyxx}$. It turns out—as you can perhaps see from the way the calculations went—that these terms are *always* equal for a cubic crystal, no matter how many force terms are taken into account, *provided* only that the forces act along the line joining each pair of atoms—that is, so long as the forces between atoms are like springs and don't have a sideways part such as you might get from a cantilevered beam (and you do get in covalent bonds).

We can check this conclusion with the experimental measurements of the elastic constants. In Table 39-2 we give the observed values of the three elastic coefficients for several cubic crystals.* You will notice that C_{xxyy} and C_{xyxy} are, in general, not equal. The reason is that in metals like sodium and potassium the interatomic forces are not along the line joining the atoms, as we assumed in our model. Diamond does not obey the law either, because the forces in diamond are covalent forces and have some directional properties—the bonds would prefer to be at the tetrahedral angle. The ionic crystals like lithium fluoride, sodium chloride, and so on, do have nearly all the physical properties assumed in our model, and the table shows that the constants C_{xxyy} and C_{xyxy} are almost equal. It is not clear why silver chloride should not satisfy the condition that $C_{xxyy} = C_{xyxy}$.

Table 39-2*
Elastic Moduli of Cubic Crystals
in 10^{12} dynes·cm²

	C_{xxxx}	C_{xxyy}	C_{xyxy}
Na	0.055	0.042	0.049
K	0.046	0.037	0.026
Fe	2.37	1.41	1.16
Diamond	10.76	1.25	5.76
Al	1.08	0.62	0.28
LiF	1.19	0.54	0.53
NaCl	0.486	0.127	0.128
KCl	0.40	0.062	0.062
NaBr	0.33	0.13	0.13
KI	0.27	0.043	0.042
AgCl	0.60	0.36	0.062

* From C. Kittel, *Introduction to Solid State Physics*, John Wiley and Sons, Inc., New York, 2nd. ed., 1956, p. 93.

* In the literature you will often find that a different notation is used. For instance, people usually write $C_{xxxx} = C_{11}$, $C_{xxyy} = C_{12}$, and $C_{xyxy} = C_{44}$.

The Flow of Dry Water

40-1 Hydrostatics

The subject of the flow of fluids, and particularly of water, fascinates everybody. We can all remember, as children, playing in the bathtub or in mud puddles with the strange stuff. As we get older, we watch streams, waterfalls, and whirlpools, and we are fascinated by this substance which seems almost alive relative to solids. The behavior of fluids is in many ways very unexpected and interesting—it is the subject of this chapter and the next. The efforts of a child trying to dam a small stream flowing in the street and his surprise at the strange way the water works its way out has its analog in our attempts over the years to understand the flow of fluids. We have tried to dam the water up—in our understanding—by getting the laws and the equations that describe the flow. We will describe these attempts in this chapter. In the next chapter, we will describe the unique way in which water has broken through the dam and escaped our attempts to understand it.

We suppose that the elementary properties of water are already known to you. The main property that distinguishes a fluid from a solid is that a fluid cannot *maintain* a shear stress for any length of time. If a shear is applied to a fluid, it will move under the shear. Thicker liquids like honey move less easily than fluids like air or water. The measure of the ease with which a fluid yields is its viscosity. In this chapter we will consider only situations in which the viscous effects can be ignored. The effects of viscosity will be taken up in the next chapter.

We begin by considering *hydrostatics*, the theory of liquids at rest. When liquids are at rest, there are no shear forces (even for viscous liquids). The law of hydrostatics, therefore, is that the stresses are always normal to any surface inside the fluid. The normal force per unit area is called the *pressure*. From the fact that there is no shear in a static fluid it follows that the pressure stress is the same in all directions (Fig. 40-1). We will let you entertain yourself by proving that if there is no shear on any plane in a fluid, the pressure must be the same in any direction.

The pressure in a fluid may vary from place to place. For example, in a static fluid at the earth's surface the pressure will vary with height because of the weight of the fluid. If the density ρ of the fluid is considered constant, and if the pressure at some arbitrary zero level is called p_0 (Fig. 40-2), then the pressure at a height h above this point is $p = p_0 - \rho gh$, where g is the gravitational force per unit mass. The combination

$$p + \rho gh$$

is, therefore, a constant in the static fluid. This relation is familiar to you, but we will now derive a more general result of which it is a special case.

If we take a small cube of water, what is the net force on it from the pressure? Since the pressure at any place is the same in all directions, there can be a net force per unit volume only because the pressure varies from one point to another. Suppose that the pressure is varying in the x -direction—and we take the coordinate directions parallel to the cube edges. The pressure on the face at x gives the force $p \Delta y \Delta z$ (Fig. 40-3), and the pressure on the face at $x + \Delta x$ gives the force $-[p + (\partial p / \partial x) \Delta x] \Delta y \Delta z$, so that the resultant force is $-(\partial p / \partial x) \Delta x \Delta y \Delta z$. If we take the remaining pairs of faces of the cube, we easily see that the pressure force per unit volume is $-\nabla p$. If there are other forces in addition—such as gravity—then the pressure must balance them to give equilibrium.

40-1 Hydrostatics

40-2 The equations of motion

40-3 Steady flow—Bernoulli's theorem

40-4 Circulation

40-5 Vortex lines

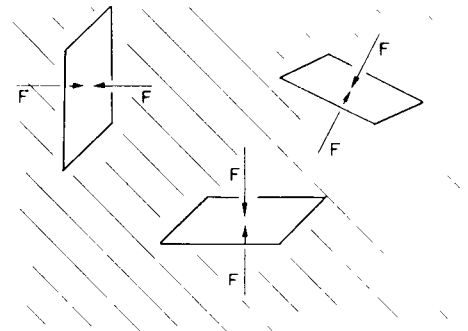


Fig. 40-1. In a static fluid the force per unit area across any surface is normal to the surface and is the same for all orientations of the surface.

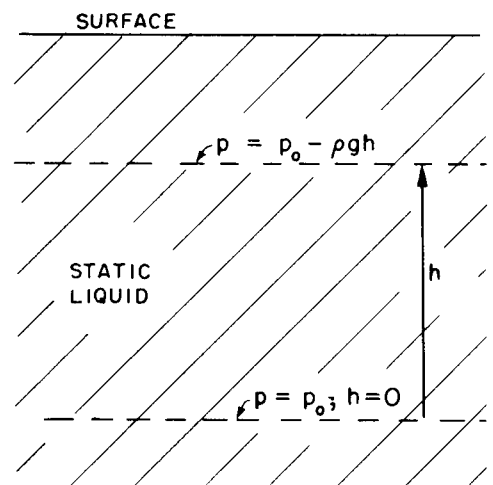


Fig. 40-2. The pressure in a static liquid.

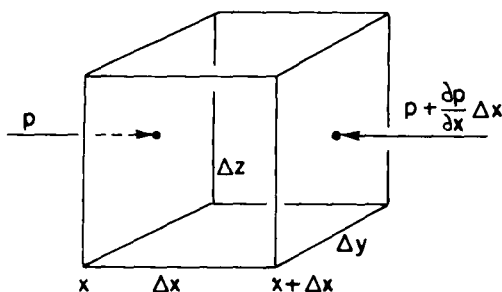


Fig. 40-3. The net pressure force on a cube is $-\nabla p$ per unit volume.

Let's take a circumstance in which such an additional force can be described by a potential energy, as would be true in the case of gravitation; we will let ϕ stand for the potential energy per unit mass. (For gravity, for instance, ϕ is just gz .) The force per unit mass is given in terms of the potential by $-\nabla\phi$, and if ρ is the density of the fluid, the force per unit volume is $-\rho\nabla\phi$. For equilibrium this force per unit volume added to the pressure force per unit volume must give zero:

$$-\nabla p - \rho\nabla\phi = 0. \quad (40.1)$$

Equation (40.1) is the equation of hydrostatics. In *general*, it *has no solution*. If the density varies in space in an arbitrary way, there is no way for the forces to be in balance, and the fluid cannot be in static equilibrium. Convection currents will start up. We can see this from the equation since the pressure term is a pure gradient, whereas for variable ρ the other term is not. Only when ρ is a constant is the potential term a pure gradient. Then the equation has a solution

$$p + \rho\phi = \text{const.}$$

Another possibility which allows hydrostatic equilibrium is for ρ to be a function only of p . However, we will leave the subject of hydrostatics because it is not nearly so interesting as the situation when fluids are in motion.

40-2 The equations of motion

First, we will discuss fluid motions in a purely abstract, theoretical way and then consider special examples. To describe the motion of a fluid, we must give its properties at every point. For example, at different places, the water (let us call the fluid "water") is moving with different *velocities*. To specify the character of the flow, therefore, we must give the three components of velocity at every point and for any time. If we can find the equations that determine the velocity, then we would know how the liquid moves at all times. The velocity, however, is not the only property that the fluid has which varies from point to point. We have just discussed the variation of the *pressure* from point to point. And there are still other variables. There may also be a variation of *density* from point to point. In addition, the fluid may be a conductor and carry an electric *current* whose density j varies from point to point in magnitude and direction. There may be a *temperature* which varies from point to point, or a *magnetic field*, and so on. So the number of fields needed to describe the complete situation will depend on how complicated the problem is. There are interesting phenomena when currents and magnetism play a dominant part in determining the behavior of the fluid; the subject is called *magnetohydrodynamics*, and great attention is being paid to it at the present time. However, we are not going to consider these more complicated situations because there are already interesting phenomena at a lower level of complexity, and even the more elementary level will be complicated enough.

We will take the situation where there is no magnetic field and no conductivity, and we will not worry about the temperature because we will suppose that the density and pressure determine in a unique manner the temperature at any point. As a matter of fact, we will reduce the complexity of our work by making the assumption that the density is a constant—we imagine that the fluid is essentially incompressible. Putting it another way, we are supposing that the variations of pressure are so small that the changes in density produced thereby are negligible. If that is not the case, we would encounter phenomena additional to the ones we will be discussing here—for example, the propagation of sound or of shock waves. We have already discussed the propagation of sound and shocks to some extent, so we will now isolate our consideration of hydrodynamics from these other phenomena by making the approximation that the density ρ is a constant. It is easy to determine when the approximation of constant ρ is a good one. We can say that if the velocities of flow are much less than the speed of a sound wave in the fluid, we do not have to worry about variations in density. The escape that water makes in our attempts to understand it is not related to the approximation of

constant density. The complications that do permit the escape will be discussed in the next chapter.

In the general theory of fluids one must begin with an *equation of state* for the fluid which connects the pressure to the density. In our approximation this equation of state is simply

$$\rho = \text{const.}$$

This then is the first relation for our variables. The next relation expresses the conservation of matter—if matter flows away from a point, there must be a decrease in the amount left behind. If the fluid velocity is \mathbf{v} , then the mass which flows in a unit time across a unit area of surface is the component of $\rho\mathbf{v}$ normal to the surface. We have had a similar relation in electricity. We also know from electricity that the divergence of such a quantity gives the rate of decrease of the density per unit time. In the same way, the equation

$$\nabla \cdot (\rho\mathbf{v}) = -\frac{\partial\rho}{\partial t} \quad (40.2)$$

expresses the conservation of mass for a fluid; it is the hydrodynamic *equation of continuity*. In our approximation, which is the incompressible fluid approximation, ρ is a constant, and the equation of continuity is simply

$$\nabla \cdot \mathbf{v} = 0. \quad (40.3)$$

The fluid velocity \mathbf{v} —like the magnetic field \mathbf{B} —has zero divergence. (The hydrodynamic equations are often closely analogous to the electrodynamic equations; that's why we studied electrodynamics first. Some people argue the other way; they think that one should study hydrodynamics first so that it will be easier to understand electricity afterwards. But electrodynamics is really much easier than hydrodynamics.)

We will get our next equation from Newton's law which tells us how the velocity changes because of the forces. The mass of an element of volume of the fluid times its acceleration must be equal to the force on the element. Taking an element of unit volume, and writing the force per unit volume as \mathbf{f} , we have

$$\rho \times (\text{acceleration}) = \mathbf{f}.$$

We will write the force density as the sum of three terms. We have already considered the pressure force per unit volume, $-\nabla p$. Then there are the “external” forces which act at a distance—like gravity or electricity. When they are conservative forces with a potential per unit mass, ϕ , they give a force density $-\rho\nabla\phi$. (If the external forces are not conservative, we would have to write \mathbf{f}_{ext} for the external force per unit volume.) Then there is another “internal” force per unit volume, which is due to the fact that in a *flowing* fluid there can also be a shearing stress. This is called the viscous force, which we will write \mathbf{f}_{visc} . Our equation of motion is

$$\rho \times (\text{acceleration}) = -\nabla p - \rho \nabla\phi + \mathbf{f}_{\text{visc}}. \quad (40.4)$$

For this chapter we are going to suppose that the liquid is “thin” in the sense that the viscosity is unimportant, so we will omit \mathbf{f}_{visc} . When we drop the viscosity term, we will be making an approximation which describes some ideal stuff rather than real water. John von Neumann was well aware of the tremendous difference between what happens when you don't have the viscous terms and when you do, and he was also aware that, during most of the development of hydrodynamics until about 1900, almost the main interest was in solving beautiful *mathematical* problems with this approximation which had almost nothing to do with real fluids. He characterized the theorist who made such analyses as a man who studied “dry water.” Such analyses leave out an *essential* property of the fluid. It is because we are leaving this property out of our calculations in this chapter that we have given it the title “The Flow of Dry Water.” We are postponing a discussion of *real* water to the next chapter.

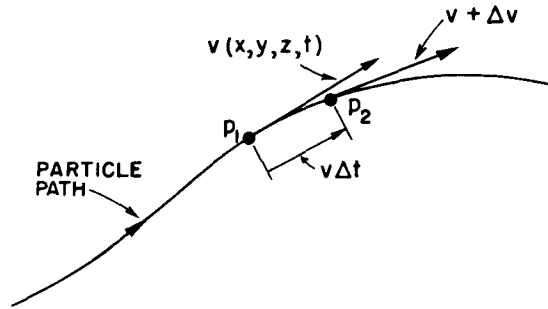


Fig. 40-4. The acceleration of a fluid particle.

If we leave out f_{visc} , we have in Eq. (40.4) everything we need except an expression for the acceleration. You might think that the formula for the acceleration of a fluid particle would be very simple, for it seems obvious that if v is the velocity of a fluid particle at some place in the fluid, the acceleration would just be $\partial v / \partial t$. *It is not*—and for a rather subtle reason. The derivative $\partial v / \partial t$, is the rate at which the velocity $v(x, y, z, t)$ changes at a *fixed point* in space. What we need is how fast the velocity changes for a *particular piece* of fluid. Imagine that we mark one of the drops of water with a colored speck so we can watch it. In a small interval of time Δt , this drop will move to a different location. If the drop is moving along some path as sketched in Fig. 40-4, it might in Δt move from P_1 to P_2 . In fact, it will move in the x -direction by an amount $v_x \Delta t$, in the y -direction by the amount $v_y \Delta t$, and in the z -direction by the amount $v_z \Delta t$. We see that, if $v(x, y, z, t)$ is the velocity of the fluid particle which is at (x, y, z) at the time t , then the velocity of the *same* particle at the time $t + \Delta t$ is given by $v(x + \Delta x, y + \Delta y, z + \Delta z, t + \Delta t)$ —with

$$\Delta x = v_x \Delta t, \quad \Delta y = v_y \Delta t, \quad \text{and} \quad \Delta z = v_z \Delta t.$$

From the definition of the partial derivatives—recall Eq. (2.7)—we have, to first order, that

$$\begin{aligned} v(x + v_x \Delta t, y + v_y \Delta t, z + v_z \Delta t, t + \Delta t) \\ = v(x, y, z, t) + \frac{\partial v}{\partial x} v_x \Delta t + \frac{\partial v}{\partial y} v_y \Delta t + \frac{\partial v}{\partial z} v_z \Delta t + \frac{\partial v}{\partial t} \Delta t. \end{aligned}$$

The acceleration $\Delta v / \Delta t$ is

$$v_x \frac{\partial v}{\partial x} + v_y \frac{\partial v}{\partial y} + v_z \frac{\partial v}{\partial z} + \frac{\partial v}{\partial t}.$$

We can write this symbolically—treating ∇ as a vector—as

$$(v \cdot \nabla)v + \frac{\partial v}{\partial t}. \quad (40.5)$$

Note that there can be an acceleration even though $\partial v / \partial t = 0$ so that velocity *at a given point* is not changing. As an example, water flowing in a circle at a constant speed is accelerating even though the velocity at a given point is not changing. The reason is, of course, that the velocity of a particular piece of water which is initially at one point on the circle has a different direction a moment later; there is a centripetal acceleration.

The rest of our theory is just mathematical—finding solutions of the equation of motion we get by putting the acceleration (40.5) into Eq. (40.4). We get

$$\frac{\partial v}{\partial t} + (v \cdot \nabla)v = -\frac{\nabla p}{\rho} - \nabla \phi, \quad (40.6)$$

where viscosity has been omitted. We can rearrange this equation by using the following identity from vector analysis:

$$(v \cdot \nabla)v = (\nabla \times v) \times v + \frac{1}{2} \nabla(v \cdot v).$$

If we now *define* a new *vector field* $\mathbf{\Omega}$, as the curl of \mathbf{v} ,

$$\mathbf{\Omega} = \nabla \times \mathbf{v}, \quad (40.7)$$

the vector identity can be written as

$$(\mathbf{v} \cdot \nabla)\mathbf{v} = \mathbf{\Omega} \times \mathbf{v} + \frac{1}{2} \nabla v^2,$$

and our equation of motion (40.6) becomes

$$\frac{\partial \mathbf{v}}{\partial t} + \mathbf{\Omega} \times \mathbf{v} + \frac{1}{2} \nabla v^2 = -\frac{\nabla p}{\rho} - \nabla \phi. \quad (40.8)$$

You can verify that Eqs. (40.6) and (40.8) are equivalent by checking that the components of the two sides of the equation are equal—and making use of (40.7).

The vector field $\mathbf{\Omega}$ is called the *vorticity*. If the vorticity is zero everywhere, we say that the flow is *irrotational*. We have already defined in Section 3–5 a thing called the *circulation* of a vector field. The circulation around any closed loop in a fluid is the line integral of the fluid velocity, at a given instant of time, around that loop:

$$(\text{Circulation}) = \oint \mathbf{v} \cdot d\mathbf{s}.$$

The circulation *per unit area* for an infinitesimal loop is then—using Stokes' theorem—equal to $\nabla \times \mathbf{v}$. So the vorticity $\mathbf{\Omega}$ is the circulation around a unit area (perpendicular to the direction of $\mathbf{\Omega}$). It also follows that if you put a little piece of dirt—*not* an infinitesimal point—at any place in the liquid it will rotate with the angular velocity $\mathbf{\Omega}/2$. Try to see if you can prove that. You can also check it out that for a bucket of water on a turntable, $\mathbf{\Omega}$ is equal to twice the local angular velocity of the water.

If we are interested only in the velocity field, we can eliminate the pressure from our equations. Taking the curl of both sides of Eq. (40.8), remembering that ρ is a constant and that the curl of any gradient is zero, and using Eq. (40.3), we get

$$\frac{\partial \mathbf{\Omega}}{\partial t} + \nabla \times (\mathbf{\Omega} \times \mathbf{v}) = 0. \quad (40.9)$$

This equation, together with the equations

$$\mathbf{\Omega} = \nabla \times \mathbf{v} \quad (40.10)$$

and

$$\nabla \cdot \mathbf{v} = 0, \quad (40.11)$$

describes completely the velocity field \mathbf{v} . Mathematically speaking, if we know $\mathbf{\Omega}$ at some time, then we know the curl of the velocity vector, and we also know that its divergence is zero, so given the physical situation we have all we need to determine \mathbf{v} everywhere. (It is just like the situation in magnetism where we had $\nabla \cdot \mathbf{B} = 0$ and $\nabla \times \mathbf{B} = \mathbf{j}/\epsilon_0 c^2$.) Thus, a given $\mathbf{\Omega}$ determines \mathbf{v} just as a given \mathbf{j} determines \mathbf{B} . Then, knowing \mathbf{v} , Eq. (40.9) tells us the rate of change of $\mathbf{\Omega}$ from which we can get the new $\mathbf{\Omega}$ for the next instant. Using Eq. (40.10), again we find the new \mathbf{v} , and so on. You see how these equations contain all the machinery for calculating the flow. Note, however, that this procedure gives the velocity field only; we have lost all information about the pressure.

We point out one special consequence of our equation. If $\mathbf{\Omega} = 0$ everywhere at any time t , $\partial \mathbf{\Omega}/\partial t$ also vanishes, so that $\mathbf{\Omega}$ is still zero everywhere at $t + \Delta t$. We have a solution to the equation; the flow is permanently irrotational. If a flow was started with zero rotation, it would always have zero rotation. The equations to be solved then are

$$\nabla \cdot \mathbf{v} = 0, \quad \nabla \times \mathbf{v} = 0.$$

They are just like the equations for the electrostatic or magnetostatic fields in free space. We will come back to them and look at some special problems later.

40-3 Steady flow—Bernoulli's theorem

Now we want to return to the equation of motion, Eq. (40.8), but limit ourselves to situations in which the flow is “steady.” By steady flow we mean that at any one place in the fluid the velocity never changes. The fluid at any point is always replaced by new fluid moving in exactly the same way. The velocity picture always looks the same— \mathbf{v} is a static vector field. In the same way that we drew “field lines” in magnetostatics, we can now draw lines which are always tangent to the fluid velocity as shown in Fig. 40-5. These lines are called *streamlines*. For steady flow, they are evidently the actual paths of fluid particles. (In unsteady flow the streamline pattern changes in time, and the streamline pattern at any instant does not represent the path of a fluid particle.)

A steady flow does not mean that nothing is happening—atoms in the fluid are moving and changing their velocities. It only means that $\partial \mathbf{v} / \partial t = 0$. Then if we take the dot product of \mathbf{v} into the equation of motion, the term $\mathbf{v} \cdot (\boldsymbol{\Omega} \times \mathbf{v})$ drops out, and we are left with

$$\mathbf{v} \cdot \nabla \left\{ \frac{p}{\rho} + \phi + \frac{1}{2} v^2 \right\} = 0. \quad (40.12)$$

This equation says that *for a small displacement in the direction of the fluid velocity* the quantity inside the brackets doesn't change. Now in steady flow all displacements are along streamlines, so Eq. (40.12) tells us that *for all the points along a streamline*, we can write

$$\frac{p}{\rho} + \frac{1}{2} v^2 + \phi = \text{const (streamline)}. \quad (40.13)$$

This is *Bernoulli's theorem*. The constant may in general be different for different streamlines; all we know is that the left-hand side of Eq. (40.13) is the same all along a *given streamline*. Incidentally, we may notice that for steady *irrotational* motion for which $\boldsymbol{\Omega} = 0$, the equation of motion (40.8) gives us the relation

$$\nabla \left\{ \frac{p}{\rho} + \frac{1}{2} v^2 + \phi \right\} = 0,$$

so that

$$\frac{p}{\rho} + \frac{1}{2} v^2 + \phi = \text{const (everywhere)}. \quad (40.14)$$

It's just like Eq. (40.13) *except* that *now* the constant has the *same value throughout the fluid*.

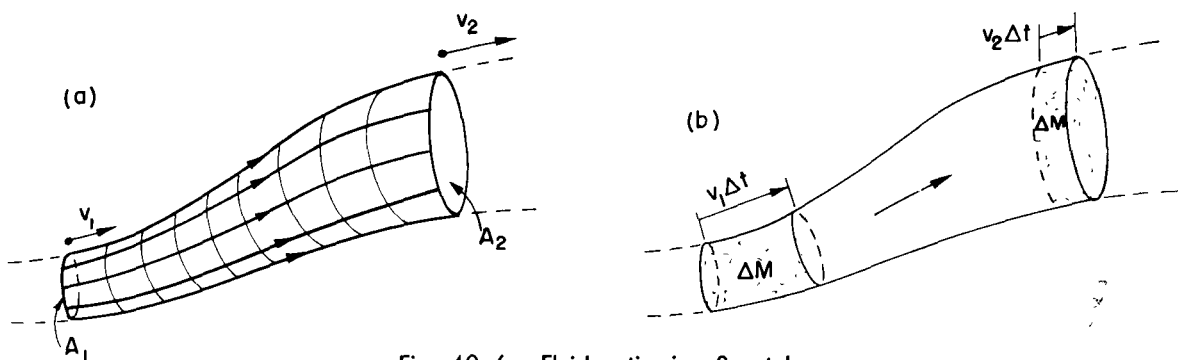


Fig. 40-6. Fluid motion in a flow tube.

The theorem of Bernoulli is in fact nothing more than a statement of the conservation of energy. A conservation theorem such as this gives us a lot of information about a flow without our actually having to solve the detailed equations. Bernoulli's theorem is so important and so simple that we would like to show you how it can be derived in a way that is different from the formal calculations we have just used. Imagine a bundle of adjacent streamlines which form a stream tube as sketched in Fig. 40-6. Since the walls of the tube consist of streamlines, no fluid flows out through the wall. Let's call the area at one end of the stream

tube A_1 , the fluid velocity there v_1 , the density of the fluid ρ_1 , and the potential energy ϕ_1 . At the other end of the tube, we have the corresponding quantities A_2 , v_2 , ρ_2 , and ϕ_2 . Now after a short interval of time Δt , the fluid at A_1 has moved a distance $v_1 \Delta t$, and the fluid at A_2 has moved a distance $v_2 \Delta t$ [Fig. 40-6(b)]. The conservation of mass requires that the mass which enters through A_1 must be equal to the mass which leaves through A_2 . These masses at these two ends must be the same:

$$\Delta M = \rho_1 A_1 v_1 \Delta t = \rho_2 A_2 v_2 \Delta t.$$

So we have the equality

$$\rho_1 A_1 v_1 = \rho_2 A_2 v_2. \quad (40.15)$$

This equation tells us that the velocity varies inversely with the area of the stream tube if ρ is constant.

Now we calculate the work done by the fluid pressure. The work done on the fluid entering at A_1 is $p_1 A_1 v_1 \Delta t$, and the work given up at A_2 is $p_2 A_2 v_2 \Delta t$. The net work on the fluid between A_1 and A_2 is, therefore,

$$p_1 A_1 v_1 \Delta t - p_2 A_2 v_2 \Delta t,$$

which must equal the increase in the energy of a mass ΔM of fluid in going from A_1 to A_2 . In other words,

$$p_1 A_1 v_1 \Delta t - p_2 A_2 v_2 \Delta t = \Delta M (E_2 - E_1), \quad (40.16)$$

where E_1 is the energy per unit mass of fluid at A_1 , and E_2 is the energy per unit mass at A_2 . The energy per unit mass of the fluid can be written as

$$E = \frac{1}{2}v^2 + \phi + U,$$

where $\frac{1}{2}v^2$ is the kinetic energy per unit mass, ϕ is the potential energy per unit mass, and U is an additional term which represents the internal energy per unit mass of fluid. The internal energy might correspond, for example, to the thermal energy in a compressible fluid, or to chemical energy. All these quantities can vary from point to point. Using this form for the energies in (40.16), we have

$$\frac{p_1 A_1 v_1 \Delta t}{\Delta M} - \frac{p_2 A_2 v_2 \Delta t}{\Delta M} = \frac{1}{2} v_2^2 + \phi_2 + U_2 - \frac{1}{2} v_1^2 - \phi_1 - U_1$$

But we have seen that $\Delta M = \rho A v \Delta t$, so we get

$$\frac{p_1}{\rho_1} + \frac{1}{2} v_1^2 + \phi_1 + U_1 = \frac{p_2}{\rho_2} + \frac{1}{2} v_2^2 + \phi_2 + U_2, \quad (40.17)$$

which is the Bernoulli result with an additional term for the internal energy. If the fluid is incompressible, the internal energy term is the same on both sides, and we get again that Eq. (40.14) holds along any streamline.

We consider now some simple examples in which the Bernoulli integral gives us a description of the flow. Suppose we have water flowing out of a hole near the bottom of a tank, as drawn in Fig. 40-7. We take a situation in which the flow speed v_{out} at the hole is much larger than the flow speed near the top of the tank; in other words, we imagine that the diameter of the tank is so large that we can neglect the drop in the liquid level. (We could make a more accurate calculation if we wished.) At the top of the tank the pressure is p_0 , the atmospheric pressure, and the pressure at the sides of the jet is also p_0 . Now we write our Bernoulli equation for a streamline, such as the one shown in the figure. At the top of the tank, we take v equal to zero and we also take the gravity potential ϕ to be zero. At the speed v_{out} , and $\phi = -gh$, so that

$$p_0 = p_0 + \frac{1}{2}\rho v_{\text{out}}^2 - \rho gh,$$

or

$$v_{\text{out}} = \sqrt{2gh}. \quad (40.18)$$

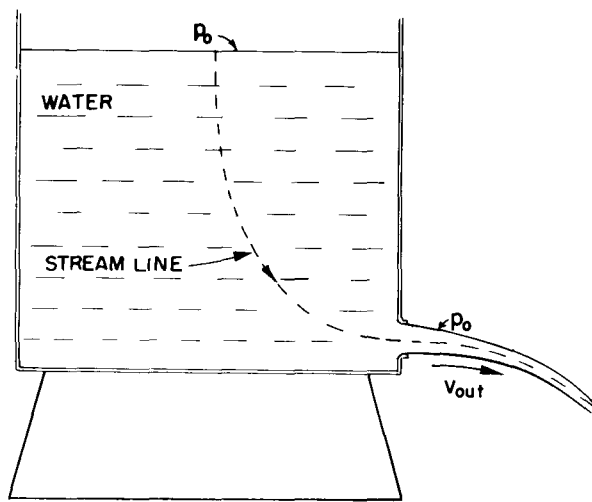


Fig. 40-7. Flow from a tank.

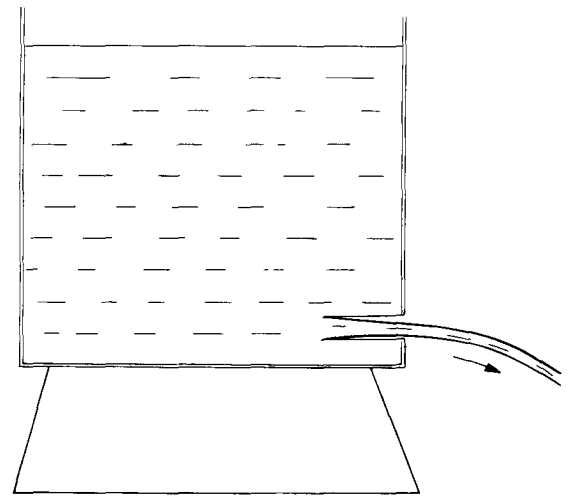


Fig. 40-8. With a re-entrant discharge tube, the stream contracts to one-half the area of the opening.

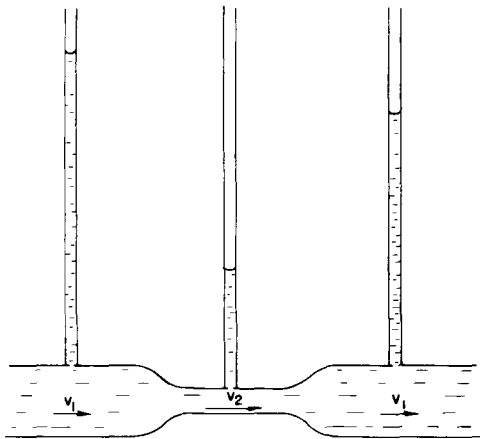


Fig. 40-9. The pressure is lowest where the velocity is highest.

This velocity is just what we would get for something which falls the distance h . It is not too surprising, since the water at the exit gains kinetic energy at the expense of the potential energy of the water at the top. Do not get the idea, however, that you can figure out the rate that the fluid flows out of the tank by multiplying this velocity by the area of the hole. The fluid velocities as the jet leaves the hole are not all parallel to each other but have components inward toward the center of the stream—the jet is converging. After the jet has gone a little way, the contraction stops and the velocities do become parallel. So the total flow is the velocity times the area *at that point*. In fact, if we have a discharge opening which is just a round hole with a sharp edge, the jet contracts to 62 percent of the area of the hole. The reduced effective area of the discharge varies for different shapes of discharge tubes, and experimental contractions are available as tables of *efflux coefficients*.

If the discharge tube is re-entrant, as shown in Fig. 40-8, it is possible to prove in a most beautiful way that the efflux coefficient is exactly 50 percent. We will give just a hint of how the proof goes. We have used the conservation of energy to get the velocity, Eq. (40.18), but there is also momentum conservation to consider. Since there is an outflow of momentum in the discharge jet, there must be a force applied over the cross section of the discharge tube. Where does the force come from? The force must come from the pressure on the walls. As long as the efflux hole is small and away from the walls, the fluid velocity near the walls of the tank will be very small. Therefore, the pressure on every face is almost exactly the same as the static pressure in a fluid at rest—from Eq. (30.14). Then the static pressure at any point on the side of the tank must be matched by an equal pressure at the point on the opposite wall, *except* at the points on the wall opposite the charge tube. If we calculate the momentum poured out through the jet by this pressure, we can show that the efflux coefficient is $1/2$. We cannot use this method for a discharge hole like that shown in Fig. 40-7, however, because the velocity increase along the wall right near the discharge area gives a pressure fall which we are not able to calculate.

Let's look at another example—a horizontal pipe with changing cross section, as shown in Fig. 40-9, with water flowing in one end and out the other. The conservation of energy, namely Bernoulli's formula, says that the pressure is lower in the constricted area where the velocity is higher. We can easily demonstrate this effect by measuring the pressure at different cross sections with small vertical columns of water attached to the flow tube through holes small enough so that they do not disturb the flow. The pressure is then measured by the height of water in these vertical columns. The pressure is found to be less at the constriction than it is on either side. If the area beyond the constriction comes back to the same value it had before the constriction, the pressure rises again.

Bernoulli's formula would predict that the pressure downstream of the constriction should be the same as it was upstream, but actually it is noticeably less. The reason that our prediction is wrong is that we have neglected the frictional, viscous forces which cause a pressure drop along the tube. Despite this pressure drop the pressure is definitely lower at the constriction (because of the increased speed) than it is on either side of it—as predicted by Bernoulli. The speed v_2 must certainly exceed v_1 to get the same amount of water through the narrower tube. So the water accelerates in going from the wide to the narrow part. The force that gives this acceleration comes from the drop in pressure.

We can check our results with another simple demonstration. Suppose we have on a tank a discharge tube which throws a jet of water upward as shown in Fig. 40-10. If the efflux velocity were exactly $\sqrt{2gh}$, the discharge water should rise to a level even with the surface of the water in the tank. Experimentally, it falls somewhat short. Our prediction is roughly right, but again viscous friction which has not been included in our energy conservation formula has resulted in a loss of energy.

Have you ever held two pieces of paper close together and tried to blow them apart? Try it! They come *together*. The reason, of course, is that the air has a higher speed going through the constricted space between the sheets than it does when it gets outside. The pressure between the sheets is *lower* than atmospheric pressure, so they come together rather than separating.

40-4 Circulation

We saw at the beginning of the last section that if we have an incompressible fluid with no circulation, the flow satisfies the following two equations:

$$\nabla \cdot \mathbf{v} = 0, \quad \nabla \times \mathbf{v} = 0. \quad (40.19)$$

They are the same as the equations of electrostatics or magnetostatics in empty space. The divergence of the electric field is zero when there are no charges, and the curl of the electrostatic field is always zero. The curl of the magnetic field is zero if there are no currents, and the divergence of the magnetic field is always zero. Therefore, Eqs. (40.19) have the same solutions as the equations for \mathbf{E} in electrostatics or for \mathbf{B} in magnetostatics. As a matter of fact, we have already solved the problem of the flow of a fluid past a sphere, as an electrostatic analogy, in Section 12-5. The electrostatic analog is a uniform electric field plus a dipole field. The dipole field is so adjusted that the flow velocity normal to the surface of the sphere is zero. The same problem for the flow past a cylinder can be worked out in a similar way by using a suitable line dipole with a uniform flow field. This solution holds for a situation in which the fluid velocity at large distances is constant—both in magnitude and direction. The solution is sketched in Fig. 40-11(a).

There is another solution for the flow around a cylinder when the conditions are such that the fluid at large distances moves in circles around the cylinder. The flow is, then, circular everywhere, as in Fig. 40-11(b). Such a flow has a circulation around the cylinder, although $\nabla \times \mathbf{v}$ is still zero *in the fluid*. How can there be circulation without a curl? We have a circulation around the cylinder because the line integral of \mathbf{v} around any loop *enclosing* the cylinder is not zero. At the same time, the line integral of \mathbf{v} around any closed path which does *not* include the cylinder is zero. We saw the same thing when we found the magnetic field around a wire. The curl of \mathbf{B} was zero outside of the wire, although a line integral of \mathbf{B} around a path which encloses the wire did not vanish. The velocity field in an irrotational circulation around a cylinder is precisely the same as the magnetic field around a wire. For a circular path with its center at the center of the cylinder, the line integral of the velocity is

$$\oint \mathbf{v} \cdot d\mathbf{s} = 2\pi r v.$$

For irrotational flow the integral must be independent of r . Let's call the constant

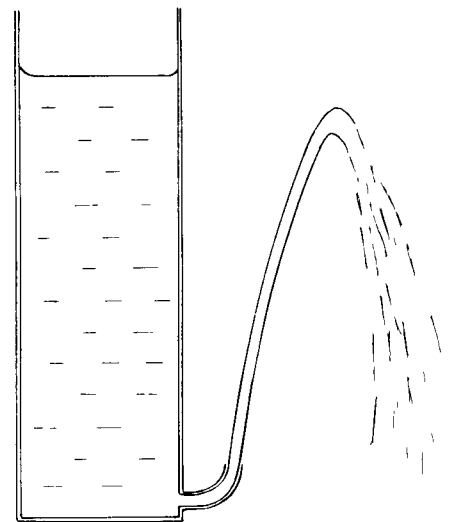


Fig. 40-10. Proof that v is not equal to $\sqrt{2gh}$.

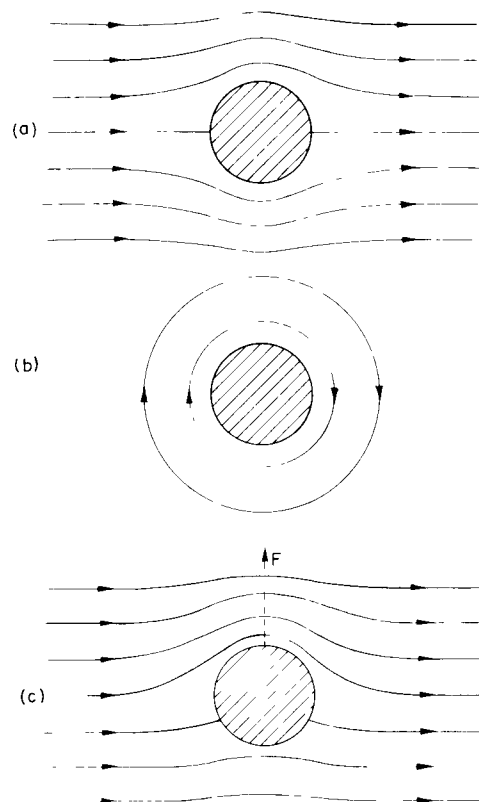


Fig. 40-11. (a) Ideal fluid flow past a cylinder. (b) Circulation around a cylinder. (c) The superposition of (a) and (b).

value C , then we have that

$$v = \frac{C}{2\pi r}, \quad (40.20)$$

where v is the tangential velocity, and r is the distance from the axis.

There is a nice demonstration of a fluid circulating around a hole. You take a transparent cylindrical tank with a drain hole in the center of the bottom. You fill it with water, stir up some circulation with a stick, and pull the drain plug. You get the pretty effect shown in Fig. 40-12. (You've seen a similar thing many times in the bathtub!) Although you put in some ω at beginning, it soon dies down because of viscosity and the flow becomes irrotational—although still with some circulation around the hole.

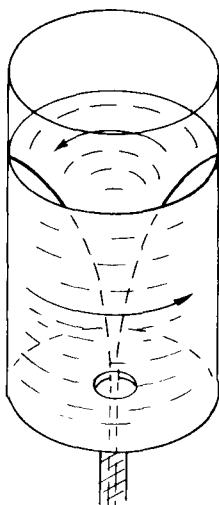


Fig. 40-12. Water with circulation draining from a tank.

From the theory, we can calculate the shape of the inner surface of the water. As a particle of the water moves inward it picks up speed. From Eq. (40.20) the tangential velocity goes as $1/r$ —it's just from the conservation of angular momentum, like the skater pulling in her arms. Also the radial velocity goes as $1/r$. Ignoring the tangential motion, we have water going radially inward toward a hole; from $\nabla \cdot \mathbf{v} = 0$, it follows that the radial velocity is proportional to $1/r$. So the total velocity also increases as $1/r$, and the water goes in along Archimedean spirals. The air-water surface is all at atmospheric pressure, so it must have—from Eq. (40.14)—the property that

$$gz + \frac{1}{2}mv^2 = \text{const.}$$

But v is proportional to $1/r$, so the shape of the surface is

$$(z - z_0) = \frac{k}{r^2}.$$

An interesting point—which is *not true in general* but is true for incompressible, irrotational flow—is that if we have one solution and a second solution, then the sum is also a solution. This is true because the equations in (40.19) are linear. The complete equations of hydrodynamics, Eqs. (40.8), (40.9), and (40.10), are not linear, which makes a vast difference. For the irrotational flow about the cylinder, however, we can superpose the flow of Fig. 40-11(a) on the flow of Fig. 40-11(b) and get the new flow pattern shown in Fig. 40-11(c). This flow is of special interest. The flow velocity is higher on the upper side of the cylinder than on the lower side. The pressures are therefore *lower* on the *upper* side than on the lower side. So when we have a combination of a circulation around a cylinder *and* a net horizontal flow, there is a net *vertical force* on the cylinder—it is called a *lift force*. Of course, if there is no circulation, there is no net force on any body according to our theory of “dry” water.

40-5 Vortex lines

We have already written down the general equations for the flow of an incompressible fluid when there may be vorticity. They are

$$\text{I. } \nabla \cdot \mathbf{v} = 0,$$

$$\text{II. } \boldsymbol{\Omega} = \nabla \times \mathbf{v},$$

$$\text{III. } \frac{\partial \boldsymbol{\Omega}}{\partial t} + \nabla \times (\boldsymbol{\Omega} \times \mathbf{v}) = 0.$$

The physical content of these equations has been described in words by Helmholtz in terms of three theorems. First, imagine that in the fluid we were to draw *vortex lines* rather than streamlines. By vortex lines we mean field lines that have the direction of $\boldsymbol{\Omega}$ and have a density in any region proportional to the magnitude of $\boldsymbol{\Omega}$. From II the divergence of $\boldsymbol{\Omega}$ is *always* zero (remember—Section 3-7—that the divergence of a curl is always zero). So vortex lines are like lines of \mathbf{B} —they never start or stop, and will tend to go in closed loops. Now Helmholtz described III

in words by the following statement: the vortex lines *move with the fluid*. This means that if you were to mark the fluid particles along some vortex lines—by coloring them with ink, for example—then as the fluid moves and carries those particles along, they will always mark the new positions of the vortex lines. In whatever way the atoms of the liquid move, the vortex lines move with them. That is one way to describe the laws.

It also suggests a method for solving any problems. Given the initial flow pattern—say \mathbf{v} everywhere—then you can calculate $\boldsymbol{\Omega}$. From the \mathbf{v} you can also tell where the vortex lines are going to be a little later—they move with the speed \mathbf{v} . With the new $\boldsymbol{\Omega}$ you can use I and II to find the new \mathbf{v} . (That's just like the problem of finding \mathbf{B} , given the currents.) If we are given the flow pattern at one instant we can in principle calculate it for all subsequent times. We have the general solution for nonviscous flow.

We would like to show how Helmholtz's statement—and, therefore, III—can be at least partly understood. It is really just the law of conservation of angular momentum applied to the fluid. Suppose we imagine a small cylinder of the liquid whose axis is parallel to the vortex lines, as in Fig. 40-13(a). At some time later, this *same* piece of fluid will be somewhere else. Generally it will occupy a cylinder with a different diameter and be in a different place. It may also have a different orientation, say as in Fig. 40-13(b). If the diameter has changed, however, the length will have increased to keep the volume constant (since we are assuming an incompressible fluid). Also, since the vortex lines are stuck with the material, their density will go up as the cross-sectional area goes down. The product of the vorticity $\boldsymbol{\Omega}$ and area A of the cylinder will remain constant, so according to Helmholtz, we should have

$$\Omega_2 A_2 = \Omega_1 A_1. \quad (40.21)$$

Now notice that with zero viscosity all the forces on the surface of the cylindrical volume (or *any* volume, for that matter) are perpendicular to the surface. The pressure forces can cause the volume to be moved from place to place, or can cause it to change shape; but with no *tangential* forces the magnitude of the *angular momentum of the material inside* cannot change. The angular momentum of the liquid in the little cylinder is its moment of inertia I times the angular velocity of the liquid, which is proportional to the vorticity $\boldsymbol{\Omega}$. For a cylinder, the moment of inertia is proportional to mr^2 . So from the conservation of angular momentum, we would conclude that

$$(M_1 R_1^2) \Omega_1 = (M_2 R_2^2) \Omega_2.$$

But the mass is the same, $M_1 = M_2$, and the areas are proportional to R^2 , so we get again just Eq. (40.21). Helmholtz's statement—which is equivalent to III—is just a consequence of the fact that in the absence of viscosity the angular momentum of an element of the fluid cannot change.

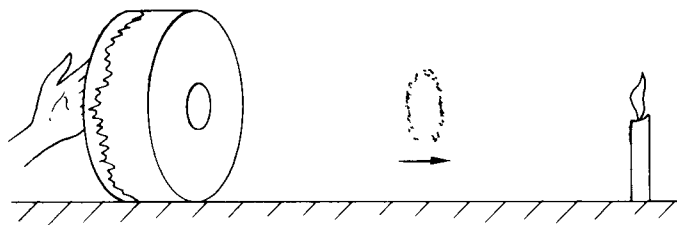


Fig. 40-14. Making a travelling vortex ring.

There is a nice demonstration of a moving vortex which is made with the simple apparatus of Fig 40-14. It is a “drum” two feet in diameter and two feet long made by stretching a thick rubber sheet over the open end of a cylindrical “box.” The “bottom”—the drum is tipped on its side—is solid except for a 3-inch diameter hole. If you give a sharp blow on the rubber diaphragm with your hand, a vortex ring is projected out of the hole. Although the vortex is invisible, you can tell it's there because it will blow out a candle 10 to 20 feet away. By the delay in

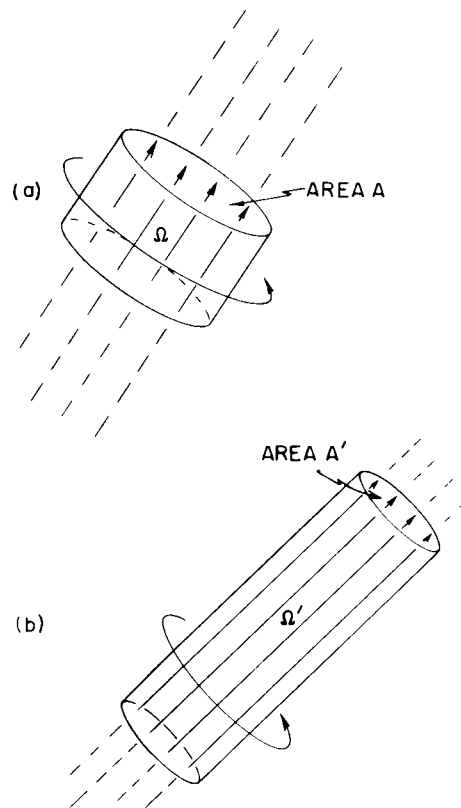


Fig. 40-13. (a) A group of vortex lines at t ; (b) the same lines at a later time t' .

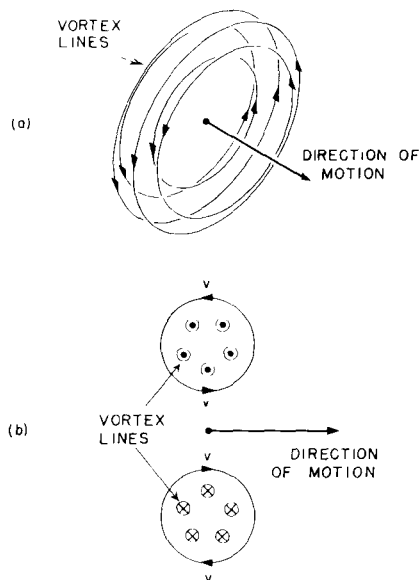


Fig. 40-15. A moving vortex ring (a smoke ring). (a) The vortex lines. (b) A cross section of the ring.

the effect, you can tell that “something” is travelling at a finite speed. You can see better what is going on if you first blow some smoke into the box. Then you see the vortex as a beautiful round “smoke ring.”

The smoke ring is a torus-shaped bundle of vortex lines, as shown in Fig. 40-15(a). Since $\Omega = \nabla \times v$, these vortex lines represent also a circulation of v as shown in part (b) of the figure. We can understand the forward motion of the ring in the following way: The circulating velocity around the *bottom* of the ring extends up to the top of the ring, having there a forward motion. Since the lines of Ω move with the fluid, they also move ahead with the velocity v . (Of course, the circulation of v around the top part of the ring is responsible for the forward motion of the vortex lines at the bottom.)

We must now mention a serious difficulty. We have already noted that Eq. (40-9) says that, if Ω is initially zero, it will always be zero. This result is a great failure of the theory of “dry” water, because it means that once Ω is zero it is *always* zero—it is impossible to *produce* any vorticity under any circumstance. Yet, in our simple demonstration with the drum, we can generate a vortex ring starting with air which was initially at rest. (Certainly, $v = 0$, $\Omega = 0$ everywhere in the box before we hit it.) Also, we all know that we can start some vorticity in a lake with a paddle. Clearly, we must go to a theory of “wet” water to get a complete understanding of the behavior of a fluid.

Another feature of the dry water theory which is incorrect is the supposition we make regarding the flow at the boundary between it and the surface of a solid. When we discussed the flow past a cylinder—as in Fig. 40-11, for example—we permitted the fluid to slide along the surface of the solid. In our theory, the velocity at a solid surface could have any value depending on how it got started, and we did not consider any “friction” between the fluid and the solid. It is an experimental fact, however, that the velocity of a real fluid always goes to zero at the surface of a solid object. Therefore, our solution for the cylinder, with or without circulation, is wrong—as is our result regarding the generation of vorticity. We will tell you about the more correct theories in the next chapter.

The Flow of Wet Water

41-1 Viscosity

In the last chapter we discussed the behavior of water, disregarding the phenomenon of viscosity. Now we would like to discuss the phenomena of the flow of fluids, *including* the effects of viscosity. We want to look at the *real behavior* of fluids. We will describe qualitatively the actual behavior of the fluids under various different circumstances so that you will get some feel for the subject. Although you will see some complicated equations and hear about some complicated things, it is not our purpose that you should learn all these things. This is, in a sense, a “cultural” chapter which will give you some idea of the way the world is. There is only one item which is worth learning, and that is the simple definition of viscosity which we will come to in a moment. The rest is only for your entertainment.

In the last chapter we found that the laws of motion of a fluid are contained in the equation

$$\frac{\partial \mathbf{v}}{\partial t} + (\mathbf{v} \cdot \nabla) \mathbf{v} = -\frac{\nabla p}{\rho} - \nabla \phi + \frac{\mathbf{f}_{\text{visc}}}{\rho}. \quad (41.1)$$

In our “dry” water approximation we left out the last term, so we were neglecting all viscous effects. Also, we sometimes made an additional approximation by considering the fluid as incompressible; then we had the additional equation

$$\nabla \cdot \mathbf{v} = 0.$$

This last approximation is often quite good—particularly when flow speeds are much slower than the speed of sound. But in real fluids it is almost never true that we can neglect the internal friction that we call viscosity; most of the interesting things that happen come from it in one way or another. For example, we saw that in “dry” water the circulation never changes—if there is none to start out with, there will never be any. Yet, circulation in fluids is an everyday occurrence. We must fix up our theory.

We begin with an important experimental fact. When we worked out the flow of “dry” water around or past a cylinder—the so-called “potential flow”—we had no reason not to permit the water to have a velocity tangent to the surface; only the normal component had to be zero. We took no account of the possibility that there might be a shear force between the liquid and the solid. It turns out—although it is not at all self-evident—that in all circumstances where it has been experimentally checked, the *velocity of a fluid is exactly zero at the surface of a solid*. You have noticed, no doubt, that the blade of a fan will collect a thin layer of dust—and that it is still there after the fan has been churning up the air. You can see the same effect even on the great fan of a wind tunnel. Why isn’t the dust blown off by the air? In spite of the fact that the fan blade is moving at high speed through the air, the speed of the air relative to the fan blade goes to zero right at the surface. So the very smallest dust particles are not disturbed.* We must modify the theory to agree with the experimental fact that in all ordinary fluids, the molecules next to a solid surface have zero velocity (relative to the surface).†

* You can blow *large* dust particles from a table top, but *not* the very finest ones. The large ones stick up into the breeze.

† You can imagine circumstances when it is not true: glass is theoretically a “liquid,” but it can certainly be made to slide along a steel surface. So our assertion must break down somewhere.

41-1 Viscosity

41-2 Viscous flow

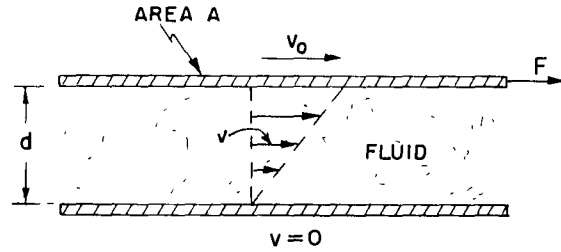
41-3 The Reynolds number

41-4 Flow past a circular cylinder

41-5 The limit of zero viscosity

41-6 Couette flow

Fig. 41-1. Viscous drag between two parallel plates.



We originally characterized a liquid by the fact that if you put a shearing stress on it—no matter how small—it would give way. It flows. In static situations, there are no shear stresses. But before equilibrium is reached—as long as you still push on it—there can be shear forces. *Viscosity* describes these shear forces which exist in a moving fluid. To get a measure of the shear forces during the motion of a fluid, we consider the following kind of experiment. Suppose that we have two solid plane surfaces with water between them, as in Fig. 41-1, and we keep one stationary while moving the other parallel to it at the slow speed v_0 . If you measure the force required to keep the upper plate moving, you find that it is proportional to the area of the plates and to v_0/d , where d is the distance between the plates. So the shear stress F/A is proportional to v_0/d :

$$\frac{F}{A} = \eta \frac{v_0}{d}.$$

The constant of proportionality η is called the *coefficient of viscosity*.

If we have a more complicated situation, we can always consider a little, flat, rectangular cell in the water with its faces parallel to the flow, as in Fig. 41-2. The shear force across this cell is given by

$$\frac{\Delta F}{\Delta A} = \eta \frac{\Delta v_x}{\Delta y} = \eta \frac{\partial v_x}{\partial y}. \quad (41.2)$$

Now, $\partial v_x / \partial y$ is the *rate of change* of the shear strain we defined in Chapter 38, so for a liquid, the shear stress is proportional to the *rate of change* of the shear strain.

In the general case we write

$$S_{xy} = \eta \left(\frac{\partial v_y}{\partial x} + \frac{\partial v_x}{\partial y} \right). \quad (41.3)$$

If there is a uniform rotation of the fluid, $\partial v_x / \partial y$ is the negative of $\partial v_y / \partial x$ and S_{xy} is zero—as it should be since there are no stresses in a uniformly rotating fluid. (We did a similar thing in defining e_{xy} in Chapter 39.) There are, of course, the corresponding expressions for S_{yz} and S_{zx} .

As an example of the application of these ideas, we consider the motion of a fluid between two coaxial cylinders. Let the inner one have the radius a and the peripheral velocity v_a , and let the outer one have radius b and velocity v_b . See Fig. 41-3. We might ask, what is the velocity distribution between the cylinders? To answer this question, we begin by finding a formula for the viscous shear in the fluid at a distance r from the axis. From the symmetry of the problem, we can assume that the flow is always tangential and that its magnitude depends only on r ; $v = v(r)$. If we watch a speck in the water at the radius r , its coordinates as a function of time are

$$x = r \cos \omega t, \quad y = r \sin \omega t,$$

where $\omega = v/r$. Then the x - and y -components of velocity are

$$v_x = -r\omega \sin \omega t = -\omega y \quad \text{and} \quad v_y = r\omega \cos \omega t = \omega x. \quad (41.4)$$

From Eq. (41.3), we have

$$S_{xy} = \eta \left[\frac{\partial}{\partial x} (x\omega) - \frac{\partial}{\partial y} (y\omega) \right] = \eta \left[x \frac{\partial \omega}{\partial x} - y \frac{\partial \omega}{\partial y} \right]. \quad (41.5)$$

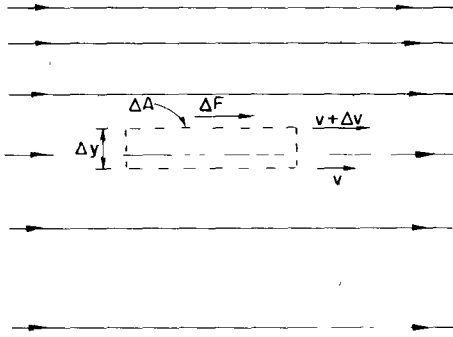


Fig. 41-2. The shear stress in a viscous fluid.

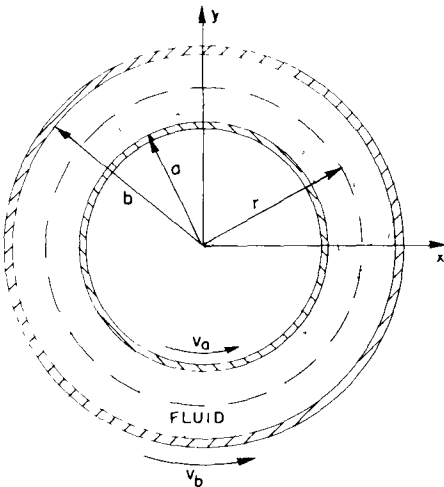


Fig. 41-3. The flow in a fluid between two concentric cylinders rotating at different angular velocities.

For a point at $y = 0$, $\partial\omega/\partial y = 0$, and $x \partial\omega/\partial x$ is the same as $r d\omega/dr$. So at that point

$$(S_{xy})_{y=0} = \eta r \frac{d\omega}{dr}. \quad (41.6)$$

(It is reasonable that S should depend on $\partial\omega/\partial r$: when there is no change in ω with r , the liquid is in uniform rotation and there are no stresses.)

The stress we have calculated is the tangential shear which is the same all around the cylinder. We can get the *torque* acting *across a cylindrical surface* at the radius r by multiplying the shear stress by the moment arm r and the area $2\pi r l$. We get

$$\tau = 2\pi r^2 l (S_{xy})_{y=0} = 2\pi \eta l r^3 \frac{d\omega}{dr}. \quad (41.7)$$

Since the motion of the water is steady—there is no angular acceleration—the net torque on the cylindrical shell of water between r and $r + dr$ must be zero; that is, the torque at r must be balanced by an equal and opposite torque at $r + dr$, so that τ must be independent of r . In other words, $r^3 d\omega/dr$ is equal to some constant, say A , and

$$\frac{d\omega}{dr} = \frac{A}{r^3}. \quad (41.8)$$

Integrating, we find that ω varies with r as

$$\omega = -\frac{A}{2r^2} + B. \quad (41.9)$$

The constants A and B are to be determined to fit the conditions that $\omega = \omega_a$ at $r = a$, and $\omega = \omega_b$ at $r = b$. We get that

$$\begin{aligned} A &= \frac{2a^2b^2}{b^2 - a^2} (\omega_b - \omega_a), \\ B &= \frac{b^2\omega_b - a^2\omega_a}{b^2 - a^2}. \end{aligned} \quad (41.10)$$

So we know ω as a function of r , and from it $v = \omega r$.

If we want the torque, we can get it from Eqs. (41.7) and (41.8):

$$\tau = 2\pi \eta l A$$

or

$$\tau = \frac{4\pi \eta l a^2 b^2}{b^2 - a^2} (\omega_b - \omega_a). \quad (41.11)$$

It is proportional to the relative angular velocities of the two cylinders. One standard apparatus for measuring the coefficients of viscosity is built this way. One cylinder—say the outer one—is on pivots but is held stationary by a spring balance which measures the torque on it, while the inner one is rotated at a constant angular velocity. The coefficient of viscosity is then determined from Eq. (41.11).

From its definition, you see that the units of η are newton-sec/m². For water at 20°C,

$$\eta = 10^3 \text{ newton-sec/m}^2.$$

It is usually more convenient to use the *specific viscosity*, which is η divided by the density ρ . The values for water and air are then comparable:

$$\begin{aligned} \text{water at } 20^\circ\text{C}, \quad \eta/\rho &= 10^{-6} \text{ m}^2/\text{sec}, \\ \text{air at } 20^\circ\text{C}, \quad \eta/\rho &= 15 \times 10^{-6} \text{ m}^2/\text{sec}. \end{aligned} \quad (41.12)$$

Viscosities usually depend strongly on temperature. For instance, for water just above the freezing point, η/ρ is 1.8 times larger than it is at 20°C.

41-2 Viscous flow

We now go to a general theory of viscous flow—at least in the most general form known to man. We already understand that the shear stress components are proportional to the spatial derivatives of the various velocity components such as $\partial v_x/\partial y$ or $\partial v_y/\partial x$. However, in the general case of a *compressible* fluid there is another term in the stress which depends on other derivatives of the velocity. The general expression is

$$S_{ij} = \eta \left(\frac{\partial v_i}{\partial x_j} + \frac{\partial v_j}{\partial x_i} \right) + \eta' \delta_{ij} (\nabla \cdot \mathbf{v}), \quad (41.13)$$

where x_i is any one of the rectangular coordinates x , y , or z , and v_i is any one of the rectangular coordinates of the velocity. (The symbol δ_{ij} is the Kronecker delta which is 1 when $i = j$ and 0 for $i \neq j$.) The additional term adds $\eta' \nabla \cdot \mathbf{v}$ to all the diagonal elements S_{ii} of the stress tensor. If the liquid is incompressible $\nabla \cdot \mathbf{v} = 0$, and this extra term doesn't appear. So it has to do with internal forces during compression. So two constants are required to describe the liquid, just as we had two constants to describe a homogeneous elastic solid. The coefficient η is the “ordinary” coefficient of viscosity which we have already encountered. It is also called the *first coefficient of viscosity* or the “shear viscosity coefficient,” and the new coefficient η' is called the *second coefficient of viscosity*.

Now we want to determine the viscous force per unit volume, f_{visc} , so we can put it into Eq. (41.1) to get the equation of motion for a real fluid. The force on a small cubical volume element of a fluid is the resultant of the forces on all the six faces. Taking them two at a time, we will get differences that depend on the derivatives of the stresses, and, therefore, on the second derivatives of the velocity. This is nice because it will get us back to a vector equation. The component of the viscous force per unit volume in the direction of the rectangular coordinate x_i is

$$\begin{aligned} (f_{\text{visc}})_i &= \sum_{j=1}^3 \frac{\partial S_{ij}}{\partial x_j} \\ &= \eta \sum_{j=1}^3 \frac{\partial}{\partial x_j} \left\{ \eta \left(\frac{\partial v_i}{\partial x_j} + \frac{\partial v_j}{\partial x_i} \right) \right\} + \frac{\partial}{\partial x_i} (\eta' \nabla \cdot \mathbf{v}) \end{aligned} \quad (41.14)$$

Usually, the variation of the viscosity coefficients with position is not significant and can be neglected. Then, the viscous force per unit volume contains only second derivatives of the velocity. We saw in Chapter 39 that the most general form of second derivatives that can occur in a vector equation is the sum of a term in the Laplacian ($\nabla \cdot \nabla \mathbf{v} = \nabla^2 \mathbf{v}$), and a term in the gradient of the divergence ($\nabla(\nabla \cdot \mathbf{v})$). Equation (41.14) is just such a sum with the coefficients η and $(\eta + \eta')$. We get

$$f_{\text{visc}} = \eta \nabla^2 \mathbf{v} + (\eta + \eta') \nabla(\nabla \cdot \mathbf{v}). \quad (41.15)$$

In the incompressible case, $\nabla \cdot \mathbf{v} = 0$, and the viscous force per unit volume is just $\eta \nabla^2 \mathbf{v}$. That is all that many people use; however, if you should want to calculate the absorption of sound in a fluid, you would need the second term.

We can now complete our general equation of motion for a real fluid. Substituting Eq. (41.15) into Eq. (41.1), we get

$$\rho \left\{ \frac{\partial \mathbf{v}}{\partial t} + (\mathbf{v} \cdot \nabla) \mathbf{v} \right\} = -\nabla p - \rho \nabla \phi + \eta \nabla^2 \mathbf{v} + (\eta + \eta') \nabla(\nabla \cdot \mathbf{v})$$

It's complicated. But that's the way nature is.

If we introduce the vorticity $\boldsymbol{\Omega} = \nabla \times \mathbf{v}$, as we did before, we can write our equation as

$$\begin{aligned} \rho \left\{ \frac{\partial \mathbf{v}}{\partial t} + \boldsymbol{\Omega} \times \mathbf{v} + \frac{1}{2} \nabla v^2 \right\} &= -\nabla p - \rho \nabla \phi + \eta \nabla^2 \mathbf{v} \\ &\quad + (\eta + \eta') \nabla(\nabla \cdot \mathbf{v}). \end{aligned} \quad (41.16)$$

We are supposing again that the only body forces acting are conservative forces like gravity. To see what the new term means, let's look at the incompressible fluid case. Then, if we take the curl of Eq. (41.16), we get

$$\frac{\partial \mathbf{\Omega}}{\partial t} + \nabla \times (\mathbf{\Omega} \times \mathbf{v}) = \frac{\eta}{\rho} \nabla^2 \mathbf{\Omega}. \quad (41.17)$$

This is like Eq. (40.9) except for the new term on the right-hand side. When the right-hand side was zero, we had the Helmholtz theorem that the vorticity stays with the fluid. Now, we have the rather complicated nonzero term on the right-hand side which, however, has straightforward physical consequences. If we disregard for the moment the term $\nabla \times (\mathbf{\Omega} \times \mathbf{v})$, we have a *diffusion equation*. The new term means that the vorticity $\mathbf{\Omega}$ *diffuses* through the fluid. If there is a large gradient in the vorticity, it will spread out into the neighboring fluid.

This is the term that causes the smoke ring to get thicker as it goes along. Also, it shows up nicely if you send a "clean" vortex (a "smokeless" ring made by the apparatus described in the last chapter) through a cloud of smoke. When it comes out of the cloud, it will have picked up some smoke, and you will see a hollow shell of a smoke ring. Some of the $\mathbf{\Omega}$ diffuses outward into the smoke, while still maintaining its forward motion with the vortex.

41-3 The Reynolds number

We will now describe the changes which are made in the character of fluid flow as a consequence of the new viscosity term. We will look at two problems in some detail. The first of these is the flow of a fluid past a cylinder—a flow which we tried to calculate in the previous chapter using the theory for nonviscous flow. It turns out that the viscous equations can be solved by man today only for a few special cases. So some of what we will tell you is based on experimental measurements—assuming that the experimental model satisfies Eq. (41.17).

The mathematical problem is this: We would like the solution for the flow of an incompressible, viscous fluid past a long cylinder of diameter D . The flow should be given by Eq. (41.17) and by

$$\mathbf{\Omega} = \nabla \times \mathbf{v} \quad (41.18)$$

with the conditions that the velocity at large distances is some constant velocity, say V (parallel to the x -axis), and at the surface of the cylinder is zero. That is,

$$v_x = v_y = v_z = 0 \quad (41.19)$$

for

$$x^2 + y^2 = \frac{D^2}{4}.$$

That specifies completely the mathematical problem.

If you look at the equations, you see that there are four different parameters to the problem: η , ρ , D , and V . You might think that we would have to give a whole series of cases for different V 's, different D 's, and so on. However, that is not the case. All the different possible solutions correspond to different values of *one parameter*. This is the most important general thing we can say about viscous flow. To see why this is so, notice first that the viscosity and density appear only in the ratio η/ρ —the *specific viscosity*. That reduces the number of independent parameters to three. Now suppose we measure all distances in the only length that appears in the problem, the diameter D of the cylinder; that is, we substitute for x , y , z , the new variables x' , y' , z' with

$$x = x'D, \quad y = y'D, \quad z = z'D.$$

Then D disappears from (41.19). In the same way, if we measure all velocities in terms of V —that is, we set $v = v'V$ —we get rid of the V , and v' is just equal to 1 at large distances. Since we have fixed our units of length and velocity, our unit

of time is now D/V ; so we should set

$$t = t' \frac{D}{V}. \quad (41.20)$$

With our new variables, the derivatives in Eq. (41.18) get changed from $\partial/\partial x$ to $(1/D) \partial/\partial x'$, and so on: so Eq. (41.18) becomes

$$\mathbf{\Omega} = \nabla \times \mathbf{v} = \frac{V}{D} \nabla' \times \mathbf{v}' = \frac{V}{D} \mathbf{\Omega}'. \quad (41.21)$$

Our main equation (41.17) then reads

$$\frac{\partial \mathbf{\Omega}'}{\partial t'} + \nabla' \times (\mathbf{\Omega}' \times \mathbf{v}') = \frac{\eta}{\rho V D} \nabla'^2 \mathbf{\Omega}'.$$

All the constants condense into one factor which we write, following tradition, as $1/\Re$:

$$\Re = \frac{\rho}{\eta} V D. \quad (41.22)$$

If we just remember that all of our equations are to be written with all quantities in the new units, we can omit all the primes. Our equations for the flow are then

$$\frac{\partial \mathbf{\Omega}}{\partial t} + \nabla \times (\mathbf{\Omega} \times \mathbf{v}) = \frac{1}{\Re} \nabla^2 \mathbf{\Omega} \quad (41.23)$$

and

$$\mathbf{\Omega} = \nabla \times \mathbf{v}$$

with the conditions

$$\mathbf{v} = 0$$

for

$$x^2 + y^2 = 1/4 \quad (41.24)$$

and

$$v_x = 1, \quad v_y = v_z = 0$$

for

$$x^2 + y^2 + z^2 \gg 1.$$

What this all means physically is very interesting. It means, for example, that if we solve the problem of the flow for one velocity V_1 and a certain cylinder diameter D_1 , and then ask about the flow for a different diameter D_2 and a different fluid, the flow will be the same for the velocity V_2 which gives the same Reynolds number—that is, when

$$\Re_1 = \frac{\rho_1}{\eta_1} V_1 D_1 = \Re_2 = \frac{\rho_2}{\eta_2} V_2 D_2. \quad (41.25)$$

For any two situations which have the same Reynolds number, the flows will “look” the same—in terms of the appropriate scaled x' , y' , z' , and t' . This is an important proposition because it means that we can determine what the behavior of the flow of air past an airplane wing will be without having to build an airplane and try it. We can, instead, make a model and make measurements using a velocity that gives the same Reynolds number. This is the principle which allows us to apply the results of “wind-tunnel” measurements on small-scale airplanes, or “model-basin” results on scale model boats, to the full-scale objects. Remember, however, that we can only do this provided the compressibility of the fluid can be neglected. Otherwise, a new quantity enters—the speed of sound. And different situations will really correspond to each other only if the ratio of V to the sound speed is also the same. This latter ratio is called the *Mach number*. So, for velocities near the speed of sound or above, the flows are the same in two situations if *both* the *Mach number* and the *Reynolds number* are the same for both situations.

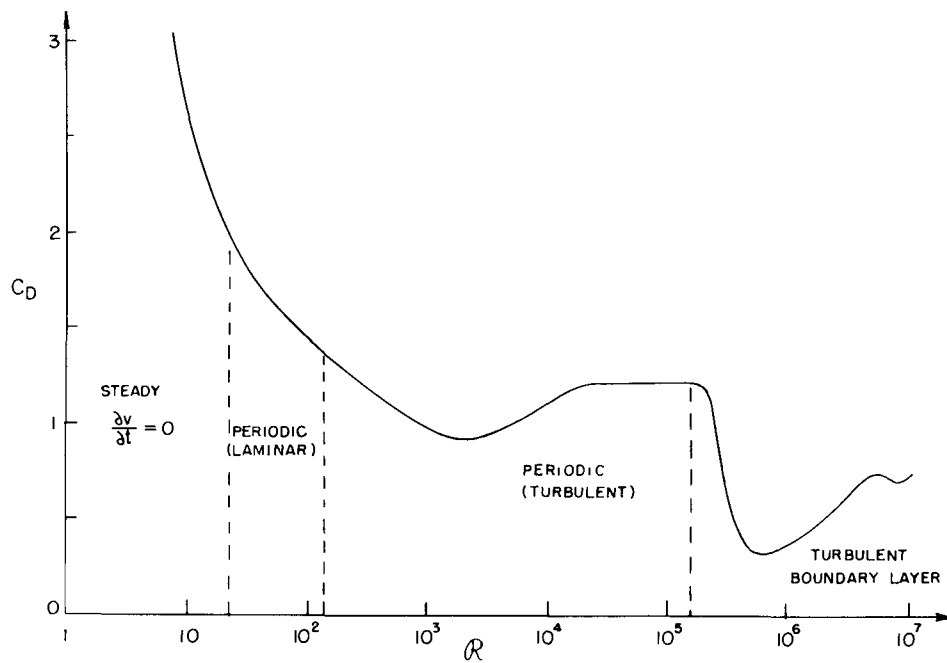


Fig. 41-4. The drag coefficient C_D of a circular cylinder as a function of the Reynolds number.

41-4 Flow past a circular cylinder

Let's go back to the problem of low-speed (nearly incompressible) flow over the cylinder. We will give a qualitative description of the flow of a real fluid. There are many things we might want to know about such a flow—for instance, what is the drag force on the cylinder? The drag force on a cylinder is plotted in Fig. 41-4 as a function of R —which is proportional to the air speed V if everything else is held fixed. What is actually plotted is the so-called *drag coefficient* C_D , which is a dimensionless number equal to the force divided by $\frac{1}{2}\rho V^2 D l$, where D is the diameter, l is the length of the cylinder, and ρ is the density of the liquid:

$$C_D = \frac{F}{\frac{1}{2}\rho V^2 D l}.$$

The coefficient of drag varies in a rather complicated way, giving us a pre-hint that something rather interesting and complicated is happening in the flow. We will now describe the nature of flow for the different ranges of the Reynolds number. First, when the Reynolds number is very small, the flow is quite steady; that is, the velocity is constant at any place, and the flow goes around the cylinder. The actual distribution of the flow lines is, however, not like it is in potential flow. They are solutions of a somewhat different equation. When the velocity is very low or, what is equivalent, when the viscosity is very high so the stuff is like honey, then the inertial terms are negligible and the flow is described by the equation

$$\nabla^2 \Omega = 0.$$

This equation was first solved by Stokes. He also solved the same problem for a sphere. If you have a small sphere moving under such conditions of low Reynolds number, the force needed to drag it is equal to $6\pi\eta aV$, where a is the radius of the sphere and V is its velocity. This is a very useful formula because it tells the speed at which tiny grains of dirt (or other particles which can be approximated as spheres) move through a fluid under a given force—as, for instance, in a centrifuge, or in sedimentation, or diffusion. In the low Reynolds number region—for R less than 1—the lines of v around a *cylinder* are as drawn in Fig. 41-5.

If we now increase the fluid speed to get a Reynolds number somewhat greater than 1, we find that the flow is different. There is a circulation behind the sphere, as shown in Fig. 41-6(b). It is still an open question as to whether there is always

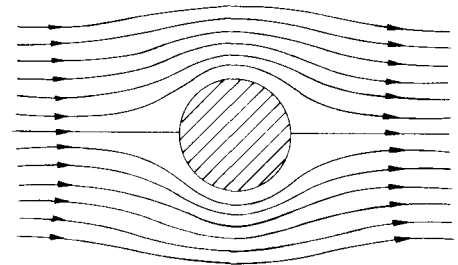


Fig. 41-5. Viscous flow (low velocities) around a circular cylinder.

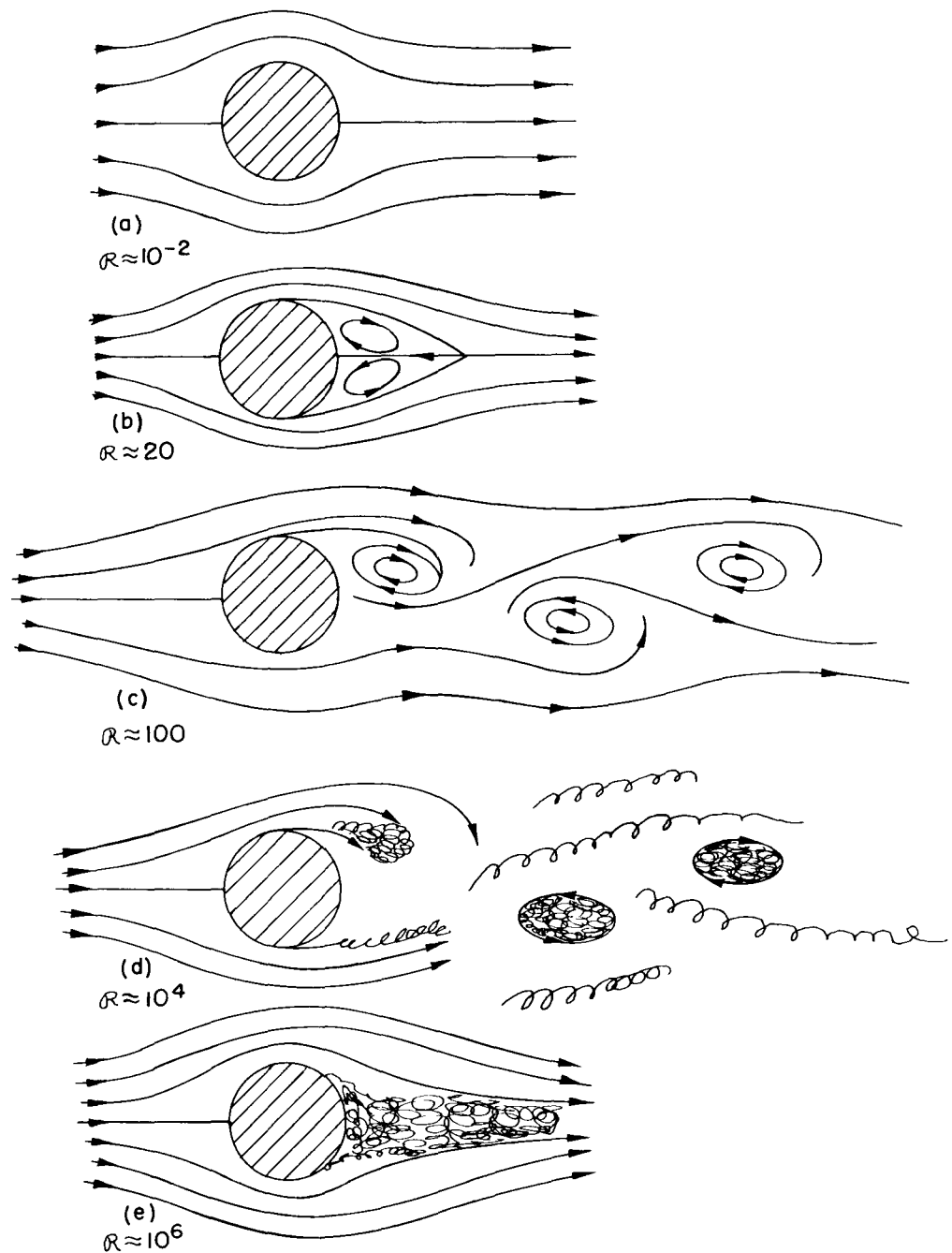


Fig. 41-6. Flow past a cylinder for various Reynolds numbers.

a circulation there even at the smallest Reynolds number or whether things suddenly change at a certain Reynolds number. It used to be thought that the circulation grew continuously. But it is now thought that it appears suddenly, and it is certain that the circulation increases with R . In any case, there is a different character to the flow for R in the region from about 10 to 30. There is a pair of vortices behind the cylinder.

The flow changes again by the time we get to a number of 40 or so. There is suddenly a complete change in the character of the motion. What happens is that one of the vortices behind the cylinder gets so long that it breaks off and travels downstream with the fluid. Then the fluid curls around behind the cylinder and makes a new vortex. The vortices peel off alternately on each side, so an instantaneous view of the flow looks roughly as sketched in Fig. 41-6(c). The stream of

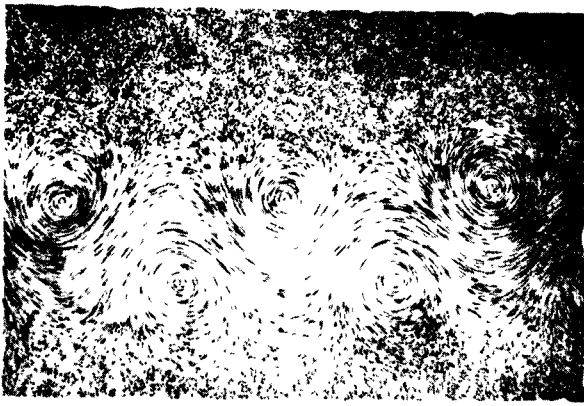


Fig. 41-7. Photograph by Ludwig Prandtl of the "vortex street" in the flow behind a cylinder.

vortices is called a "Kármán vortex street." They always appear for $\mathcal{R} > 40$. We show a photograph of such a flow in Fig. 41-7.

The difference between the two flows in Fig. 41-6(c) and 41-6(b) or 41-6(a) is almost a complete difference in regime. In Fig. 41-6(a) or (b), the velocity is constant, whereas in Fig. 41-6(c), the velocity at any point varies with time. There is no steady solution above $\mathcal{R} = 40$ —which we have marked on Fig. 41-4 by a dashed line. For these higher Reynolds numbers, the flow varies with time but in a *regular*, cyclic fashion.

We can get a physical idea of how these vortices are produced. We know that the fluid velocity must be zero at the surface of the cylinder and that it also increases rapidly away from that surface. Vorticity is created by this large local variation in fluid velocity. Now when the main stream velocity is low enough, there is sufficient time for this vorticity to diffuse out of the thin region near the solid surface where it is produced and to grow into a large region of vorticity. This physical picture should help to prepare us for the next change in the nature of the flow as the main stream velocity, or \mathcal{R} , is increased still more.

As the velocity gets higher and higher, there is less and less time for the vorticity to diffuse into a larger region of fluid. By the time we reach a Reynolds number of several hundred, the vorticity begins to fill in a thin band, as shown in Fig. 41-6(d). In this layer the flow is chaotic and irregular. The region is called the *boundary layer* and this irregular flow region works its way farther and farther upstream as \mathcal{R} is increased. In the turbulent region, the velocities are very irregular and "noisy"; also the flow is no longer two-dimensional but twists and turns in all three dimensions. There is still a regular alternating motion superimposed on the turbulent one.

As the Reynolds number is increased further, the turbulent region works its way forward until it reaches the point where the flow lines leave the cylinder—for flows somewhat above $\mathcal{R} = 10^5$. The flow is as shown in Fig. 41-6(e), and we have what is called a "turbulent boundary layer." Also, there is a drastic change in the drag force; it drops by a large factor, as shown in Fig. 41-4. In this speed region, the drag force actually *decreases* with increasing speed. There seems to be little evidence of periodicity.

What happens for still larger Reynolds numbers? As we increase the speed further, the wake increases in size again and the drag increases. The latest experiments—which go up to $\mathcal{R} = 10^7$ or so—indicate that a new periodicity appears in the wake, either because the whole wake is oscillating back and forth in a gross motion or because some new kind of vortex is occurring together with an irregular noisy motion. The details are as yet not entirely clear, and are still being studied experimentally.

41-5 The limit of zero viscosity

We would like to point out that none of the flows we have described are anything like the potential flow solution we found in the preceding chapter. This is, at first sight, quite surprising. After all, \mathcal{R} is proportional to $1/\eta$. So η going to zero is equivalent to \mathcal{R} going to infinity. And if we take the limit of large \mathcal{R} in

Eq. (41.23), we get rid of the right-hand side and get just the equations of the last chapter. Yet, you would find it hard to believe that the highly turbulent flow at $\mathcal{R} = 10^7$ was approaching the smooth flow computed from the equations of “dry” water. How can it be that as we approach $\mathcal{R} = \infty$, the flow described by Eq. (41.23) gives a completely different solution from the one we obtained taking $\eta = 0$ to start out with? The answer is very interesting. Note that the right-hand term of Eq. (41.23) has $1/\mathcal{R}$ times a *second derivative*. It is a higher derivative than any other derivative in the equation. What happens is that although the coefficient $1/\mathcal{R}$ is small, there are very rapid variations of Ω in the space near the surface. These rapid variations compensate for the small coefficient, and the product *does not go to zero* with increasing \mathcal{R} . The solutions do not approach the limiting case as the coefficient of $\nabla^2 \Omega$ goes to zero.

You may be wondering, “What is the fine-grain turbulence and how does it maintain itself? How can the vorticity which is made somewhere at the edge of the cylinder generate so much noise in the background?” The answer is again interesting. Vorticity has a tendency to amplify itself. If we forget for a moment about the diffusion of vorticity which causes a loss, the laws of flow say (as we have seen) that the vortex lines are carried along with the fluid, at the velocity v . We can imagine a certain number of lines of Ω which are being distorted and twisted by the complicated flow pattern of v . This pulls the lines closer together and mixes them all up. Lines that were simple before will get knotted and pulled close together. They will be longer and tighter together. The strength of the vorticity will increase and its irregularities—the pluses and minuses—will, in general, increase. So the magnitude of vorticity in three dimensions increases as we twist the fluid about.

You might well ask, “When is the potential flow a satisfactory theory at all?” In the first place, it is satisfactory outside the turbulent region where the vorticity has not entered appreciably by diffusion. By making special streamlined bodies, we can keep the turbulent region as small as possible; the flow around airplane wings—which are carefully designed—is almost entirely true potential flow.

41-6 Couette flow

It is possible to demonstrate that the complex and shifting character of the flow past a cylinder is not special but that the great variety of flow possibilities occurs generally. We have worked out in Section 1 a solution for the viscous flow between two cylinders, and we can compare the results with what actually happens. If we take two concentric cylinders with an oil in the space between them and put a fine aluminum powder as a suspension in the oil, the flow is easy to see. Now if we turn the outer cylinder slowly, nothing unexpected happens; see Fig. 41-8(a). Alternatively, if we turn the inner cylinder slowly, nothing very striking occurs. However, if we turn the inner cylinder at a higher rate, we get a surprise. The fluid breaks into horizontal bands, as indicated in Fig. 41-8(b). When the outer cylinder rotates at a similar rate with the inner one at rest, no such effect occurs. How can it be that there is a difference between rotating the inner or the out cylinder? After all, the flow pattern we derived in Section 1 depended only on $\omega_b - \omega_a$. We can get the answer by looking at the cross sections shown in Fig. 41-9. When the inner layers of the fluid are moving more rapidly than the outer ones, they tend to move *outward*—the centrifugal force is larger than the pressure holding them in place. A whole layer cannot move out uniformly because the outer layers are in the way. It must break into cells and circulate, as shown in Fig. 41-9(b). It is like the convection currents in a room which has hot air at the bottom. When the inner cylinder is at rest and the outer cylinder has a high velocity, the centrifugal forces build up a pressure gradient which keeps everything in equilibrium—see Fig. 41-9(c) (as in a room with hot air at the top).

Now let's speed up the inner cylinder. At first, the number of bands increases. Then suddenly you see the bands become wavy, as in Fig. 41-8(c), and the waves travel around the cylinder. The speed of these waves is easily measured. For high rotation speeds they approach $1/3$ the speed of the inner cylinder. And no one

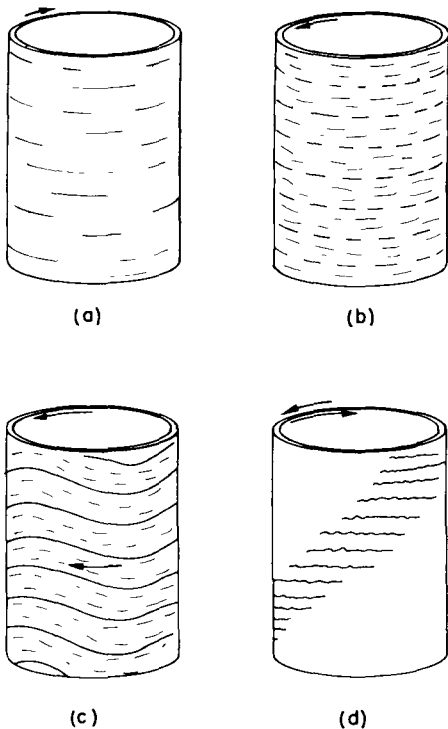


Fig. 41-8. Liquid flow patterns between two transparent rotating cylinders.

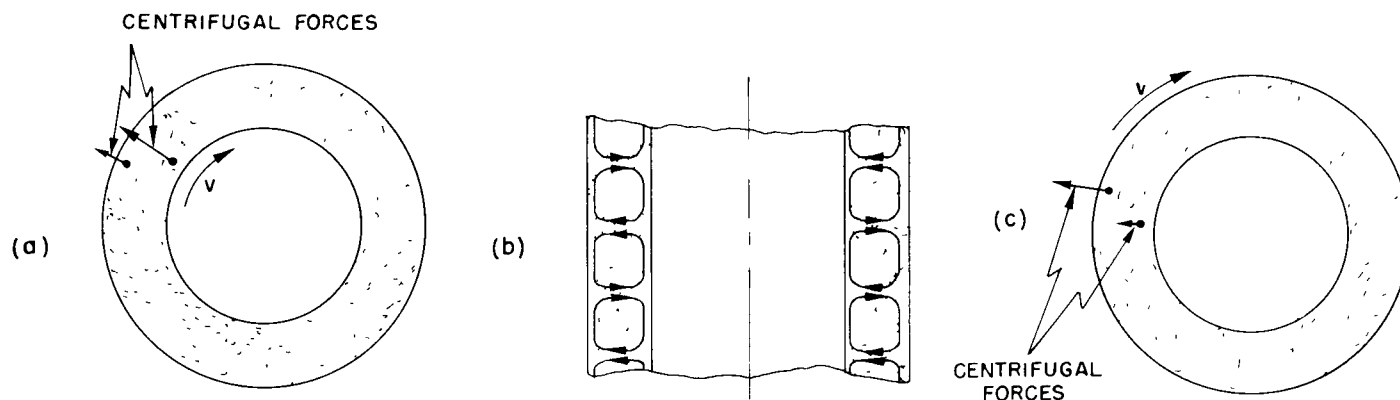


Fig. 41-9. Why the flow breaks up into bands.

knows why! There's a challenge. A simple number like $1/3$, and no explanation. In fact, the whole mechanism of the wave formation is not very well understood, yet it is steady laminar flow.

If we now start rotating the outer cylinder also—but in the opposite direction—the flow pattern starts to break up. We get wavy regions alternating with apparently quiet regions, as sketched in Fig. 41-8(d), making a spiral pattern. In these “quiet” regions, however, we can see that the flow is really quite irregular; it is, in fact completely turbulent. The wavy regions also begin to show irregular turbulent flow. If the cylinders are rotated still more rapidly, the whole flow becomes chaotically turbulent.

In this simple experiment we see many interesting regimes of flow which are quite different, and yet which are all contained in our simple equation for various values of the one parameter \mathcal{R} . With our rotating cylinders, we can see many of the effects which occur in the flow past a cylinder: first, there is a steady flow, second, a flow sets in which varies in time but in a regular, smooth way; finally, the flow becomes completely irregular. You have all seen the same effects in the column of smoke rising from a cigarette in quiet air. There is a smooth steady column followed by a series of twistings as the stream of smoke begins to break up, ending finally in an irregular churning cloud of smoke.

The main lesson to be learned from all of this is that a tremendous variety of behavior is hidden in the simple set of equations in (41.23). All the solutions are for the same equations, only with different values of \mathcal{R} . We have no reason to think that there are any terms missing from these equations. The only difficulty is that we do not have the mathematical power today to analyze them except for very small Reynolds numbers—that is, in the completely viscous case. That we have written an equation does not remove from the flow of fluids its charm or mystery or its surprise.

If such variety is possible in a simple equation with only one parameter, how much more is possible with more complex equations! Perhaps the fundamental equation that describes the swirling nebulae and the condensing, revolving, and exploding stars and galaxies is just a simple equation for the hydrodynamic behavior of nearly pure hydrogen gas. Often, people in some unjustified fear of physics say you can't write an equation for life. Well, perhaps we can. As a matter of fact, we very possibly already have the equation to a sufficient approximation when we write the equation of quantum mechanics:

$$H\psi = -\frac{\hbar}{i} \frac{\partial \psi}{\partial t}.$$

We have just seen that the complexities of things can so easily and dramatically escape the simplicity of the equations which describe them. Unaware of the scope of simple equations, man has often concluded that nothing short of God, not mere equations, is required to explain the complexities of the world.

We have written the equations of water flow. From experiment, we find a set of concepts and approximations to use to discuss the solution—vortex streets, turbulent wakes, boundary layers. When we have similar equations in a less familiar situation, and one for which we cannot yet experiment, we try to solve the equations in a primitive, halting, and confused way to try to determine what new qualitative features may come out, or what new qualitative forms are a consequence of the equations. Our equations for the sun, for example, as a ball of hydrogen gas, describe a sun without sunspots, without the rice-grain structure of the surface, without prominences, without coronas. Yet, all of these are really in the equations; we just haven't found the way to get them out.

There are those who are going to be disappointed when no life is found on other planets. Not I—I want to be reminded and delighted and surprised once again, through interplanetary exploration, with the infinite variety and novelty of phenomena that can be generated from such simple principles. The test of science is its ability to predict. Had you never visited the earth, could you predict the thunderstorms, the volcanos, the ocean waves, the auroras, and the colorful sunset? A salutary lesson it will be when we learn of all that goes on on each of those dead planets—those eight or ten balls, each agglomerated from the same dust cloud and each obeying exactly the same laws of physics.

The next great era of awakening of human intellect may well produce a method of understanding the *qualitative* content of equations. Today we cannot. Today we cannot see that the water flow equations contain such things as the barber pole structure of turbulence that one sees between rotating cylinders. Today we cannot see whether Schrödinger's equation contains frogs, musical composers, or morality—or whether it does not. We cannot say whether something beyond it like God is needed, or not. And so we can all hold strong opinions either way.

Index

- Aberration, I-27-7, I-34-10
Absolute zero, I-1-5
Absorption, I-31-8 ff
Absorption coefficient, II-32-8
Acceleration, I-8-8 ff
 components of, I-9-3
 of gravity, I-9-4
Accelerator guide field, II-29-4 ff
Activation energy, I-42-7
Active circuit element, II-22-5
Adams, J. C., I-7-5
Adiabatic compression, I-39-5
Adiabatic demagnetization, II-35-9 f
Adiabatic expansion, I-44-5
Affective future, I-17-4
Aharanov, II-15-12
Air trough, I-10-5
Algebra, I-22-1 ff
Alternating-current circuits, II-22-1 ff
Alternating-current generator,
 II-17-6 ff
Alnico V, II-37-10
Amber, II-1-10
Ammeter, II-16-1
Ampere, A., II-13-3
Ampere's law, II-13-4
Amperian current, II-36-2
Amplitudes of oscillation, I-21-3
Amplitude modulation, I-48-3
Analog computer, I-25-8
Anderson, C. D., I-52-10
Angle, of incidence, I-26-3
 of precession, II-34-4
 of reflection, I-26-3
Angstrom (unit), I-1-3
Angular frequency, I-21-3, I-29-2
Angular momentum, I-7-7, I-18-5 f,
 I-20-1
 conservation of, I-4-7, I-18-6 ff,
 I-20-5
 of rigid body, I-20-8
Anomalous refraction, I-33-9 f
Antiferromagnetic material, II-37-11
Antimatter, I-52-10 f
Antiparticle, I-2-8
Aristotle, I-5-1
Atom, I-1-2
 metastable, I-42-10
 Rutherford-Bohr model, II-5-3
 stability of, II-5-3
 Thompson model, II-5-3
Atomic clock, I-5-5
Atomic currents, II-13-5 f
Atomic hypothesis, I-1-2
Atomic orbits, II-1-8
Atomic particles, I-2-9 f
Atomic polarizability, II-32-2
Atomic processes, I-1-5 f
Attenuation, I-31-8
Avogadro, A., I-39-2
Avogadro's number, I-41-10
Axial vector, I-52-6 f
Barkhausen effect, II-37-9
Battery, II-22-6
Becquerel, A. H., I-28-3
Bell, A. G., II-16-3
Benzene molecule, III-49-10 ff
Bernoulli's theorem, II-40-6 ff
Bessel function, II-23-6
Betatron, II-17-5
Biot-Savart law, II-14-10
Birefringence, I-33-3 ff
Blackbody radiation, I-41-5 f
Boehm, I-52-10
Bohm, II-7-7, II-15-12
Bohr, N., I-42-9, II-5-3
Bohr magneton, II-34-12
Bohr radius, I-38-6
Boltzmann, L., I-41-2
Boltzmann's law, I-40-2 f
Bopp, II-28-8
Born, M., I-37-1, I-38-9, II-28-7
Boundary layer, II-41-9
Boundary-value problems, II-7-1
Boyle's law, I-40-8
"Boys" camera, II-9-10
Bragg, L., II-30-9
Bragg-Nye crystal model, II-30-9 ff
Breaking-drop theory, II-9-9
Bremsstrahlung, I-34-6 f
Brewster's angle, I-33-6
Briggs, H., I-22-6
Brown, R., I-41-1
Brownian motion, I-1-8, I-6-5,
 I-41-1 ff
Brush discharge, II-9-9
Bulk modulus, II-38-3
Calculus, differential, I-8-4, II-2-1 ff
 integral, II-3-1 ff
 of variations, II-19-3
Cantilever beam, II-38-10
Capacitance, I-23-5
 mutual, II-22-17
Capacitor, I-14-9, I-23-5, II-22-3 ff,
 II-23-2 ff
 parallel-plate, I-14-9, II-6-11 ff,
 II-8-3
Capacity, II-6-12
 of a condenser, II-8-2
Capillary action, I-51-8
Carnot, S., I-4-2, I-44-2 ff
Carnot cycle, I-44-5 f, I-45-2
Carrier signal, I-48-3
Catalyst, I-42-8
Cavendish, H., I-7-9
Cavendish's experiment, I-7-9
Cavity resonator, II-23-1 ff
Center of mass, I-18-1 f, I-19-1 ff
Centrifugal force, I-7-5, I-12-11
Čerenkov, P. A., I-51-2
Čerenkov radiation, I-51-2
Charge, conservation of, I-4-7,
 II-13-1 f
 on electron, I-12-7
 line of, II-5-3 f
 motion of, II-29-1 ff
 sheet of, II-5-4
 sphere of, II-5-4 f
Charge density, II-5-4
Charge separation, II-9-7 ff
Charged conductor, II-8-2 ff
Chemical energy, I-4-2
Chemical kinetics, I-42-7 f
Chemical reaction, I-1-6 ff
Chromaticity, I-35-6 f
Circuits, alternating-current, II-22-1 ff
 equivalent, II-22-10 f
Circuit elements, II-23-1 f
 active, II-22-5
 passive, II-22-5
Circular motion, I-21-4
Circulation, II-1-5, II-3-8 ff
Classical electron radius, II-28-3
Clausius, R., I-44-2, I-44-3
Clausius-Clapeyron equation, I-45-6 ff
Clausius-Mossotti equation, II-11-6 f,
 II-32-7
Cleavage plane, II-30-1
Coaxial line, II-24-1
Coefficient, absorption, II-32-8
 of coupling, II-17-14
 of friction, I-12-4
 gravitational, I-7-9
 of viscosity, II-41-2
Collision, I-16-6
 elastic, I-10-7
Colloidal particles, II-7-8 ff
Color vision, I-35-1 ff
 physiochemistry of, I-35-9 f
Complex impedance, I-23-7
Complex numbers, I-22-7 ff, I-23-1 ff
Complex variable, II-7-2 ff
Compound eye, I-36-6 ff
Compression, adiabatic, I-39-5
 isothermal, I-44-5
Condenser, parallel-plate, I-14-9,
 II-6-11 ff, II-8-3
Conductivity, II-32-10
 thermal, II-2-8, II-12-2
Conductor, II-1-2
Cones, I-35-1
Conservation, of angular momentum,
 I-4-7, I-18-6 ff, I-20-5
 of charge, I-4-7, II-13-1 f
 of energy, I-3-2, I-4-1 ff, II-27-1 f
 of linear momentum, I-4-7,
 I-10-1 ff

Contraction hypothesis, I-15-3
 Copernicus, I-7-1
 Coriolis force, I-19-8 f
 Cornea, I-35-1
 Cosmic rays, II-9-2
 Couette flow, II-41-10 ff
 Coulomb's law, I-28-2, II-4-2 ff,
 II-5-6
 Coupling, coefficient of, II-17-14
 Covalent bond, II-30-2
 Cross product, II-2-8, II-31-8
 Cross section for scattering, I-32-7
 Crystal, II-30-1 ff
 geometry of, II-30-1 f
 Crystal diffraction, I-38-4 f
 Crystal lattice, II-30-3 f
 Cubic cell, II-30-7
 Curie law, II-11-5
 Curie temperature, II-36-13
 Curie-Weiss law, II-11-9
 Curl operator, II-2-8, II-3-1
 Current, Amperian, II-36-2
 atomic, II-13-5 f
 eddy, II-16-6
 electric, II-13-1 f
 induced, II-16-1 ff
 Current density, II-13-1
 Cutoff frequency, II-22-14

D'Alembertian, II-25-8
 Debye length, II-7-9
 Dedekind, R., I-22-4
 Degrees of freedom, I-25-2, I-39-12
 Demagnetization, adiabatic, II-35-9 f
 Density, I-1-4
 Derivative, I-8-5 ff
 partial, I-14-9
 Diamagnetism, II-34-1 ff
 Dicke, R. H., I-7-11
 Dielectric, II-10-1 ff, II-11-1 ff
 Dielectric constant, II-10-1 f
 Differential calculus, I-8-4, II-2-1 f
 Diffraction, I-30-1 ff
 by screen, I-31-10 f
 Diffraction grating, I-29-5, I-30-3 ff
 Diffusion, I-43-1 ff
 of neutrons, II-12-6 ff
 Dipole, II-21-5 ff
 electric, II-6-2 ff
 magnetic, II-14-7 f
 Dipole moment, I-12-6, II-6-7
 Dipole potential, II-6-4 ff
 Dipole radiator, I-28-5 f, I-29-3 ff
 Dirac, P., I-52-10, II-2-1, II-28-7
 Dirac equation, I-20-6
 Dislocation, II-30-8, II-30-9
 Dispersion, I-31-6 ff
 Distance, I-5-5 ff
 Distance measurement, color brightness,
 I-5-6
 triangulation, I-5-6
 Divergence, II-25-7
 Divergence operator, II-2-7, II-3-1
 Domain, II-37-6
 Doppler effect, I-17-8, I-23-9,
 I-34-7 f, I-38-6
 Dot product, II-2-4, II-25-3
 Double stars, I-7-6
 Drag coefficient, II-41-7
 "Dry" water, II-40-1 ff

INDEX 2

Dynamics, I-7-2 f, I-9-1 ff
 relativistic, I-15-9 f

Eddy current, II-16-6
 Efficiency of ideal engine, I-44-7 f
 Einstein, A., I-2-6, I-7-11, I-12-12,
 I-15-1, I-16-1, I-41-8, I-42-8,
 I-42-9
 Elastic collision, I-10-7
 Elastic constants, II-39-6, II-39-10 f
 Elastic energy, I-4-2, I-4-6
 Elastic materials, II-39-1 ff
 Elastica, II-38-12
 Elasticity, II-38-1 ff
 Elasticity tensor, II-39-4 ff
 Electret, II-11-8
 Electric charge density, II-2-8, II-4-3
 Electric current, II-13-1 f
 in the atmosphere, II-9-2 f
 Electric current density, II-2-8
 Electric dipole, II-6-2 ff
 Electric field, I-2-4, I-12-7 f, II-1-2,
 II-1-3, II-6-1 ff, II-7-1 ff
 relativity of, II-13-6 ff
 Electric flux, II-1-4
 Electric potential, II-4-4
 Electric susceptibility, II-10-4
 Electrical energy, I-4-2, II-15-3 ff
 Electrical forces, II-1-1 ff, II-13-1
 Electrodynamics, II-1-3
 relativistic notation, II-25-1 ff
 Electromagnet, II-36-9 ff
 Electromagnetic energy, I-29-2
 Electromagnetic field, I-2-2, I-2-5,
 I-10-9
 Electromagnetic mass, II-28-3 f
 Electromagnetic radiation, I-26-1,
 I-28-1 ff
 Electromagnetic waves, II-21-1 f
 cosmic rays, I-2-5
 gamma rays, I-2-5
 infrared, I-2-5, I-23-8, I-26-1
 light, I-2-5
 ultraviolet, I-2-5, I-26-1
 x-rays, I-2-5, I-26-1
 Electromagnetism, II-1-1 ff
 laws of, II-1-5 ff
 Electromotive force, II-16-2
 Electron, I-2-4, I-37-1, I-37-4 ff
 charge on, I-12-7
 radius of, classical, I-32-4
 Electron cloud, I-6-11
 Electron microscope, II-29-3 f
 Electron-ray tube, I-12-9
 Electron volt (unit), I-34-4
 Electronic polarization, II-11-1 ff
 Electrostatic energy, II-8-1 ff
 of charges, II-8-1 f
 of ionic crystal, II-8-4 ff
 in nuclei, II-8-6 ff
 of a point charge, II-8-12
 Electrostatic equations, II-10-6 f
 Electrostatic field, II-5-1 ff, II-7-1 f
 energy in, II-8-9 ff
 of a grid, II-7-10 f
 Electrostatic lens, II-29-2 f
 Electrostatic potential, equations of,
 II-6-1
 Electrostatics, II-4-1 ff, II-5-1
 Ellipse, I-7-1

Emissivity, II-6-14
 Energy, II-22-11 f
 chemical, I-4-2
 of a condenser, II-8-2 ff
 conservation of, I-3-2, I-4-1 ff,
 II-27-1 f
 elastic, I-4-2, I-4-6
 electrical, I-4-2, II-15-3 ff
 electromagnetic, I-29-2
 electrostatic, II-8-1 ff
 in electrostatic field, II-8-9 ff
 gravitational, I-4-2 ff
 heat, I-4-2, I-4-6, I-10-7, I-10-8
 kinetic, I-1-7, I-4-2, I-4-5 f,
 I-39-4
 magnetic, II-17-12 ff
 mass, I-4-2, I-4-7
 mechanical, II-15-3 ff
 nuclear, I-4-2
 potential, I-4-4, I-13-1 ff, I-14-1 ff
 radiant, I-4-2
 relativistic, I-16-1 ff
 Energy density, II-27-2
 Energy flux, II-27-2
 Energy levels, I-38-7 f
 Energy theorem, I-50-7 f
 Enthalpy, I-45-5
 Entropy, I-44-10 ff, I-46-7 ff
 Eotvös, L., I-7-11
 Equilibrium, I-1-6
 Equipotential surfaces, II-4-11 f
 Equivalent circuits, II-22-10 f
 Euclid, I-5-6
 Euclidean geometry, I-12-3
 Euler force, II-38-11
 Evaporation, I-1-5 f
 of a liquid, I-40-3 f, I-42-1 ff
 Exchange force, II-37-2
 Excited state, II-8-7
 Expansion, adiabatic, I-44-5
 isothermal, I-44-5
 Exponential atmosphere, I-40-1 f
 Eye, compound, I-36-6 ff
 human, I-35-1 f, I-36-3 ff

Farad (unit), I-25-7, II-6-13
 Faraday, M., II-10-1
 Faraday's law of induction, II-17-2
 Fermat, P., I-26-3
 Fermi (unit), I-5-10
 Fermi, E., I-5-10
 Ferrite, II-37-12
 Ferroelectricity, II-11-8 ff
 Ferromagnetic insulators, II-37-12
 Ferromagnetism, II-34-1 f, II-36-1 ff,
 II-37-1 ff
 Feynman, R., II-28-8
 Fields, I-2-2, I-2-4, I-2-5, I-10-9,
 I-12-7 ff, I-13-8 f, I-14-7 ff
 in a cavity, II-5-8 f
 of a charged conductor, II-6-8
 of a conductor, II-5-7 f
 electric, I-2-4, I-12-7 f, II-1-2,
 II-1-3, II-6-1 ff, II-7-1 ff
 electrostatic, II-5-1 ff, II-7-1 f
 magnetic, II-1-2, II-1-3, II-13-1,
 II-14-1 ff
 magnetizing, II-36-7
 scalar, II-2-2 ff
 superposition of, I-12-9

- two-dimensional, II-7-2 ff
- vector, II-1-4 f, II-2-1 ff
- Field energy, II-27-1 ff
 - of a point charge, II-28-1 f
- Field index, II-29-5
- Field-ion microscope, II-6-14
- Field lines, II-4-11
- Field momentum, II-27-9 ff
 - of a moving charge, II-28-2 f
- Field strength, II-1-4
- Filter, II-22-14 ff
- Flow, fluid, II-12-8 ff
 - irrotationa^l, II-40-5
 - viscous, II-41-4 f
- Fluid flow, II-12-8 ff
- Flux, II-4-7 ff
 - electric, II-1-4
 - of a vector field, II-3-2 ff
- Flux rule, II-17-1 ff
- Focal length, I-27-1 ff
- Focus, I-26-5
- Force, centrifugal, I-7-5, I-12-11
 - components of, I-9-3
 - conservative, I-14-3 ff
 - Coriolis, I-19-8 f
 - electrical, I-2-3 ff, II-1-1 ff, I-13-1
 - electromotive, II-16-2
 - gravitational, I-2-3
 - Lorentz, II-13-1, II-15-14
 - magnetic, II-1-2, II-13-1
 - molecular, I-1-3, I-12-6 f
 - moment of, I-18-5
 - nonconservative, I-14-6 f
 - nuclear, I-12-12
 - pseudo, I-12-10 ff
- Fourier, J., I-50-2 f
- Fourier analysis, I-50-2 ff
- Fourier theorem, II-7-11
- Fourier transform, I-25-4
- Four-vectors, I-15-8 f, I-17-5 ff, II-25-1 ff
- Fovea, I-35-1
- Frank, I., I-51-2
- Franklin, B., II-5-6
- Frequency, angular, I-21-3, I-29-2
 - of oscillation, I-2-5
 - plasma, II-7-6, II-32-12
- Fresnel's reflection formulas, I-33-8
- Friction, I-10-5, I-12-3 ff
 - coefficient of, I-12-4

- Galileo, I-5-1, I-7-2, I-9-1, I-52-3
- Galilean relativity, I-10-3
- Galilean transformation, I-12-11
- Galvanometer, II-1-8, II-16-1
- Garnet, II-37-12
- Gauss (unit), I-34-4
- Gauss, K., II-16-2
- Gauss' law, II-4-9 f, II-5-1 ff
- Gauss' theorem, II-3-5
- Gaussian surface, II-10-1
- Geiger, II-5-3
- Gell-Mann, M., I-2-9
- Generator, alternating-current, II-17-6 ff
 - electric, II-16-1 ff, II-22-5 ff
 - van de Graaff, II-5-9, II-8-7
- Geometrical optics, I-26-1, I-27-1 f
- Gerlach, II-35-3
- Gradient operator, II-2-4, II-3-1

- Gravitation, I-2-3, I-7-1 ff, I-12-2
- Gravitational acceleration, I-9-4
- Gravitational coefficient, I-7-9
- Gravitational energy, I-4-2 ff
- Gravitational field, I-12-8 ff, I-13-8 f
- Gravity, I-13-3 ff
 - acceleration of, I-9-4
- Green's function, I-25-4
- Ground state, II-8-7
- Gyroscope, I-20-5 ff

- Hamilton's first principal function, II-19-8
- Harmonic motion, I-21-4, I-23-1 ff
- Harmonic oscillator, I-10-1, I-21-1 ff
 - forced, I-21-5 f, I-23-3 ff
- Harmonics, I-50-1 ff
- Heat, I-1-3, I-13-3
- Heat conduction, II-3-6 ff
- Heat diffusion equation, II-3-8
- Heat energy, I-4-2, I-4-6, I-10-7, I-10-8
- Heat engines, I-44-1 ff
- Heat flow, II-2-8 f, II-12-2 ff
- Heisenberg, W., I-6-10, I-37-1, I-37-9, I-37-11, I-37-12, I-38-9
- Helmholtz, H., I-35-7, II-40-10
- Henry (unit), I-25-7
- Hess, II-9-2
- Hexagonal cell, II-30-7
- High-voltage breakdown, II-6-13 f
- Hooke's law, I-12-6, II-38-1 f
- Huygens, C., I-15-2, I-26-2
- Hydrodynamics, II-40-2 ff
- Hydrostatics, II-40-1 ff
- Hypocycloid, I-34-3
- Hysteresis curve, II-37-5 ff
- Hysteresis loop, II-36-8

- Ideal gas law, I-39-10 ff
- Illumination, II-12-10 ff
- Image charge, II-6-9
- Impedance, I-25-8 f, II-22-1 ff
 - complex, I-23-7
- Incidence, angle of, I-26-3
- Inclined plane, I-4-4
- Index of refraction, I-31-1 ff
- Induced currents, II-16-1 ff
- Inductance, I-23-6, II-16-4 f, II-17-12 ff, II-22-2 f
 - mutual, II-17-9 ff, II-22-16
 - self-, II-16-4, II-17-11 f
- Induction, laws of, II-17-1 ff
- Inductor, I-23-6
- Inertia, I-2-3, I-7-11
 - moment of, I-18-7, I-19-5 ff
 - principle of, I-9-1
- Infeld, II-28-7
- Infrared radiation, I-23-8, I-26-1
- Integral, I-8-7 f
- Integral calculus, II-3-1 ff
- Insulator, II-1-2, II-10-1
- Interference, I-28-6, I-29-1 ff
- Interfering waves, I-37-4
- Interferometer, I-15-5
- Internal reflection, II-33-12
- Ion, I-1-6
- Ionic bond, II-30-2
- Ionic conductivity, I-43-6 f
- Ionic polarizability, II-11-8

- Ionization energy, I-42-5
- Ionosphere, II-7-5, II-9-3
- Irrotational flow, II-40-5
- Isotherm, II-2-3
- Isothermal atmosphere, I-40-2
- Isothermal compression, I-44-5
- Isothermal expansion, I-44-5
- Isothermal surfaces, II-2-3
- Isotopes, I-3-4 ff

- Jeans, J., I-40-9, I-41-6 f, II-2-6
- Johnson noise, I-41-2, I-41-8
- Joule (unit), I-13-3
- Joule heating, I-24-2

- Kármán vortex street, II-41-9
- Kepler, J., I-7-1
- Kepler's laws, I-7-1 f, I-9-1, I-18-6
- Kerr cell, I-33-5
- Kilocalorie (unit), II-8-5
- Kinetic energy, I-1-7, I-4-2, I-4-5 f, I-39-4
 - rotational, I-19-7 ff
- Kinetic theory, I-42-1 ff
 - of gases, I-39-1 ff
- Kirchhoff's laws, I-25-9, II-22-7 ff
- Kronecker delta, II-31-6

- Lamb, II-5-6
- Lamé elastic constants, II-39-6
- Landé *g*-factor, II-34-4
- Laplace, P., I-47-7
- Laplace equation, II-6-1, II-7-1
- Laplacian operator, II-2-10
- Larmor frequency, II-34-7
- Larmor's theorem, II-34-6 f
- Laser, I-32-6, I-42-10
- Laughton, II-5-6
- Laws, of electromagnetism, II-1-5 ff
 - of induction, II-17-1 ff
- Least action, principle of, II-19-1 ff
- Least time, principle of, I-26-3 ff, I-26-8
- Leibnitz, G. W., I-8-4
- Lens formula, I-27-6
- Lenz's rule, II-16-4, II-34-2
- Leverrier, U., I-7-5
- Liénard-Wiechert potentials, II-21-11
- Light, II-21-1 f
 - momentum of, I-34-10 f
 - polarized, I-32-9
 - scattering of, I-32-5 ff
 - speed of, I-15-1, II-18-8 f
- Light waves, I-48-1
- Lightning, II-9-10 f
- Line of charge, II-5-3 f
- Line integral, II-3-1
- Linear momentum, conservation of, I-4-7, I-10-1 ff
- Linear systems, I-25-1 ff
- Lodestone, II-1-10
- Logarithms, I-22-4
- Lorentz, H. A., I-15-3
- Lorentz condition, II-25-9
- Lorentz contraction, I-15-7
- Lorentz force, II-13-1, II-15-14
- Lorentz formula, II-21-12 f
- Lorentz gauge, II-18-11
- Lorentz transformation, I-15-3, I-17-1, I-34-8, -52-2, II-25-1
 - of fields, II-26-1 ff

- McCullough, II-1-9
 Mach number, II-41-6
 Magnetic dipole, II-14-7 f
 Magnetic dipole moment, II-14-8
 Magnetic energy, II-17-12 ff
 Magnetic field I-12-9 f, II-1-2, II-1-3, II-13-1, II-14-1 ff
 relativity of, II-13-6 ff
 of steady currents, II-13-3 f
 Magnetic force, II-1-2, II-13-1
 on a current, II-13-2 f
 Magnetic induction, I-12-10
 Magnetic lens, II-29-3
 Magnetic materials, II-37-1 ff
 Magnetic moments, II-34-3 f
 Magnetic resonance, II-35-1 ff
 Magnetic susceptibility, II-35-7
 Magnetism, I-2-4, II-34-1 ff
 Magnetization currents, II-36-1 ff
 Magnetizing fields, II-36-7
 Magnetostatics, II-4-1, II-13-1 ff
 Magnetostriction, II-37-6
 Magnification, I-27-5
 Marsden, II-5-3
 Maser, I-42-10
 Mass, I-9-1, I-15-1
 center of, I-18-1 f, I-19-1 ff
 electromagnetic, II-28-3 f
 relativistic, I-16-6 ff
 Mass energy, I-4-2, I-4-7
 Mass-energy equivalence, I-15-10 f
 Maxwell, J. C., I-6-1, I-6-9, I-28-1, I-40-8, I-41-7, I-46-5, II-1-8, II-1-11, II-5-6, II-18-1 ff
 Maxwell's equations, I-15-2, I-25-3, I-47-7, II-2-1, II-2-8, II-4-1, II-6-1, II-18-1 ff, II-32-3 ff
 currents and charges, II-21-1 ff
 free space, II-20-1 ff
 Mayer, J. R., I-3-2
 Mean free path, I-43-3 f
 Mean square distance, I-6-5, I-41-9
 Mechanical energy, II-15-3 ff
 Mendeléev, I-2-9
 Metastable atom, I-42-10
 Meter (unit), I-5-10
 Mev (unit), I-2-9
 Michelson-Morley experiment, I-15-3 ff
 Miller, W. C., I-35-2
 Minkowski, I-17-8
 Minkowski space, II-31-12
 Modes, I-49-1 ff
 Mossbauer, R., I-23-9
 Mole (unit), I-39-10
 Molecular attraction, I-1-3, I-12-6 f
 Molecular crystal, II-30-2
 Molecular diffusion, I-43-7 ff
 Molecular dipole, II-11-1
 Molecular motion, I-41-1
 Molecule, I-1-3
 Moment, dipole, I-12-6
 of force, I-18-5
 of inertia, I-18-7, I-19-5 ff
 Momentum, I-9-1 f, I-38-2 ff
 angular, I-7-7, I-18-5 ff, I-20-1, I-20-5
 of light, I-34-10 f
 linear, I-4-7, I-10-1 ff
 relativistic, I-10-8 f, I-16-1 ff
 Momentum spectrometer, II-29-1
 Momentum spectrum, II-29-2
 Monatomic gas, I-39-5
 Monoclinic cell, II-30-7
 Motion, I-5-1, I-8-1 ff
 of charge, II-29-1 ff
 circular, I-21-4
 constrained, I-14-3
 harmonic, I-21-4, I-23-1 ff
 parabolic, I-8-10
 planetary, I-7-1 ff, I-9-6 f, I-13-5
 Motors, electric, II-16-1 ff
 Moving charge, field momentum of, II-28-2 f
 Music, I-50-1
 Mutual inductance, II-17-9 ff, II-22-16

 Nernst heat theorem, I-44-11
 Neuman, J. von, II-12-9
 Neutrons, I-2-4
 diffusion of, II-12-6 ff
 Neutron diffusion equation, II-12-7
 Newton, I., I-8-4, I-15-1, I-37-1, II-4-10
 Newton-meter (unit), I-13-3
 Newton's laws, I-2-6, I-7-3 ff, I-7-11, I-9-1 ff, I-10-1 ff, I-11-7 f, I-12-1, I-39-2, I-41-1, I-46-1, II-7-5
 Nishijima, I-2-9
 Nodes, I-49-2
 Noise, I-50-1
 Nonpolar molecule, II-11-1
 Nuclear cross section, I-5-9
 Nuclear energy, I-4-2
 Nuclear forces, I-12-12
 Nuclear *g*-factor, II-34-4
 Nuclear interactions, II-8-7
 Nuclear magnetic resonance, II-35-10 ff
 Nucleus, I-2-4, I-2-8 ff
 Numerical analysis, I-9-6
 Nutation, I-20-7
 Nye, J. F., II-30-9

 Oersted (unit), II-36-6
 Ohm (unit), I-25-7
 Ohm's law, I-25-7, I-43-7
 Operator, curl, II-2-8, II-3-1
 divergence, II-2-7, II-3-1
 gradient, II-2-4, II-3-1
 Laplacian, II-2-10
 vector, II-2-6
 Optic axis, I-33-3
 Optic nerve, I-35-2
 Optics, I-26-1 ff
 geometrical, I-26-1, I-27-1 ff
 Orbital motion, II-34-3
 Orientation polarization, II-11-3 ff
 Oriented magnetic moment, II-35-4
 Orthorhombic cell, II-30-7
 Oscillation, amplitude, of, I-21-3
 damped, I-24-3 f
 frequency of, I-2-5
 period of, I-21-3
 periodic, I-9-4
 phase of, I-21-3

 Oscillator, I-5-2
 harmonic, I-10-1, I-21-1, I-21-5 f, I-23-3 ff

 Pappus, theorem of I-19-4
 Parabolic antenna, I-30-6 f
 Parabolic motion, I-8-10
 Parallel-axis theorem, I-19-6
 Parallel-plate capacitor, I-14-9, II-6-11 ff, II-8-3
 Paramagnetism, II-34-1 ff, II-35-1 ff
 Paraxial rays, I-27-2
 Partial derivative, I-14-9
 Particles, "strange", II-8-7
 Permalloy, II-37-11
 Permeability, II-36-9
 Pascal's triangle, I-6-4
 Passive circuit element, II-22-5
 Pendulum, I-49-6 f
 Pendulum clock, I-5-2
 Period of oscillation, I-21-3
 Periodic time, I-5-1 f
 Perpetual motion, I-46-2
 Phase of oscillation, I-21-3
 Phase shift, I-21-3
 Phase velocity, I-48-6
 Photon, I-2-7, I-26-1, I-37-8
 Physicochemistry of color vision, I-35-9 f
 Piezoelectricity, II-11-8
 Pines, II-7-7
 Planck, M., I-41-6, I-42-8, I-42-9
 Planck's constant, I-5-10, I-6-10, I-17-8, I-37-11
 Plane lattice, II-30-5
 Plane waves, II-21-1 ff
 Planetary motion, I-7-1 ff, I-9-6 f, I-13-5
 Plasma frequency, II-7-6, II-32-12
 Plasma oscillations, II-7-5 ff
 Plimpton, II-5-6
 Poincaré, H., I-15-3, I-15-5, I-16-1
 Poincaré stress, II-28-4
 Point charge, electrostatic energy of, II-8-12
 field energy of, II-28-1 f
 Poisson's ratio, II-38-2
 Polar molecule, II-11-1, II-11-3 ff
 Polarization, I-33-1 ff, II-32-1 ff
 Polarization charges, II-10-3 ff
 Polarization vector, II-10-2 f
 Polarized light, I-32-9
 Potential energy, I-4-4, I-13-1 ff, I-14-1 ff
 Potential gradient of the atmosphere, II-9-2 f
 Power, I-13-2
 Poynting, J., II-27-3
 Precession, angle of, II-34-4
 of atomic magnets, II-34-4 f
 Pressure, I-1-3
 Priestly, J., II-5-6
 Principle of least action, II-19-1 ff
 Principle of superposition, II-1-3, II-4-2
 Probability, I-6-1 ff
 Probability density, I-6-8 f
 Probability distribution, I-6-7 ff
 Propagation factor, II-22-14
 Proton, I-2-4

Proton spin, II-8-7
 Pseudo force, I-12-10 ff
 Ptolemy, I-26-2
 Purkinje effect, I-35-2
 Pyroelectricity, II-11-8

 Quadrupole lens, II-7-4, II-29-6
 Quadrupole potential, II-6-8
 Quantized magnetic states, II-35-1 ff
 Quantum electrodynamics, I-2-7, I-28-3
 Quantum mechanics, I-2-2, I-2-6 ff, I-6-10, I-10-9, I-37-1 ff, I-38-1 ff

 Rabi, I I., II-35-4
 Rabi molecular-beam method, II-35-4 ff
 Radiant energy, I-4-2
 Radiation, infrared, I-23-8, I-26-1 relativistic effects, I-35-1 ff synchrotron, I-34-3 ff, I-34-6 ultraviolet, I-26-1
 Radiation damping, I-32-3 f
 Radiation resistance, I-32-1 ff
 Radioactive clock I-5-3 ff
 Radius of electron, I-32-4
 Ramsey, N., I-5-5
 Random walk, I-6-5 ff, I-41-8 ff
 Ratchet and pawl machine, I-46-1 ff
 Rayleigh's criterion, I-30-6
 Rayleigh's law, I-41-6
 Rayleigh waves, II-38-8
 Reactance, II-22-11
 Reciprocity principle, I-30-7
 Rectification, I-50-9
 Rectifier, II-22-15
 Reflected waves, II-33-7 ff
 Reflection, I-26-2 f angle of, I-26-3 internal, II-33-12 of light, II-33-1 ff
 Refraction, I-26-2 f anomalous, I-33-9 f index of, I-31-1 ff of light, II-33-1 ff
 Refractive index, II-32-1 ff
 Relative permeability, II-36-9
 Relativistic dynamics, I-15-9 f
 Relativistic energy, I-16-1 ff
 Relativistic mass, I-16-1 ff
 Relativistic momentum, I-10-8 f, I-16-1 ff
 Relativity, of electric field, II-13-6 ff Galilean, I-10-3 of magnetic field, II-13-6 ff special theory of, I-15-1 ff theory of, I-7-11, I-17-1
 Resistance, I-23-5
 Resistor, I-23-5, II-22-4
 Resonant cavity, II-23-6 ff
 Resonant circuits, II-23-10 f
 Resonant mode, II-23-10
 Resonator, cavity, II-23-1 ff
 Resolving power, I-27-7 f, I-30-5 f
 Resonance, I-23-1 ff electrical, I-23-5 ff in nature, I-23-7 ff
 Resonance interaction, I-2-9
 Retarded time, I-28-2
 Retherford, II-5-6
 Retina, I-35-1
 Reynolds' number, II-41-5 f
 Rigid body, I-18-1 angular momentum of, I-20-8 rotation of, I-18-2 ff
 Ritz combination principle, I-38-8
 Rods, I-35-1, I-36-6
 Roemer, O., I-7-5
 Root-mean-square distance, I-6-6
 Rotation, of axes, I-11-3 f plane, I-18-1 of a rigid body, I-18-2 ff in space, I-20-1 ff in two dimensions, I-18-1 ff
 Rushton, I-35-9
 Rutherford, II-5-3
 Rutherford-Bohr atomic model, II-5-3
 Rydberg (unit), I-38-6

 Scalar, I-11-5
 Scalar field, II-2-2 ff
 Scalar product, II-25-3 ff
 Scattering of light, I-32-5 ff
 Schrodinger, E., I-35-6, I-37-1, I-38-9
 Schrödinger equation, II-15-12
 Scientific method, I-2-1 f
 Screw dislocation, II-30-9
 Screw jack, I-4-5
 Second (unit), I-5-5
 Seismograph, I-51-5
 Self-inductance, II-16-4, II-17-11 f
 Shannon, C , I-44-2
 Shear modulus, II-38-5
 Shear wave, I-51-4, II-38-8
 Sheet of charge, II-5-4
 Side bands, I-48-4 f
 Simultaneity, I-15-7 f
 Sinusoidal waves, I-29-2 f
 Skin depth, II-32-11
 Slip dislocation, II-30-9
 Smoluchowski, I-41-8
 Smooth muscle, I-14-2
 Snell, W., I-26-3
 Snell's law, I-26-3, I-31-2, II-33-1
 Solenoid, II-13-5
 Solid-state physics, II-8-6
 Sound, I-2-3, I-47-1 ff, I-50-1 speed of, I-47-7 f
 Space, I-8-2
 Space-time, I-2-6, I-17-1 ff, II-26-12
 Special theory of relativity, I-15-1 ff
 Specific heat, I-40-7 f, I-45-2, II-37-4
 Speed, I-8-2 ff, I-9-2 of light, I-15-1, II-18-8 f of sound, I-47-7 f
 Sphere of charge, II-5-4 f
 Spherical waves, II-20-12 ff, II-21-2 ff
 Spinel, II-37-12
 Spin orbit, II-8-7
 Spontaneous emission, I-42-9
 Standard deviation, I-6-9
 Statics, II-4-1 f
 Statistical fluctuations, I-6-3 ff
 Statistical mechanics, I-3-1, I-40-1 ff
 Steady flow, II-40-6 ff
 Step leader, II-9-10
 Stern, II-35-3
 Stern-Gerlach experiment, II-35-3 ff
 Stevinus, S., I-4-5
 Stokes' theorem, II-3-10
 Strain, II-38-2
 Strain tensor, II-31-11, II-39-1 ff
 "Strange" particles, II-8-7
 "Strangeness" number, I-2-9
 Streamlines, II-40-6
 Stress, II-38-2
 Stress tensor, II-31-9 ff
 Striated muscle, I-14-2
 Supermalloy, II-36-9
 Superposition, II-13-11 f of fields, I-12-9 principle of, I-25-2 ff, II-1-3, II-4-2
 Surface, equipotential, II-4-11 f gaussian, II-10-1 isothermal, II-2-3
 Surface tension, II-12-5
 Symmetry, I-1-4, I-11-1 ff of physical laws, I-16-3, I-52-1 ff
 Synchrotron, I-2-5, I-15-9, I-34-3 ff, I-34-6, II-17-5

 Tamm, I., I-51-2
 Taylor expansion, II-6-7
 Temperature, I-39-6 ff
 Tensor, II-26-7, II-31-1 ff
 Tensor field, II-31-11
 Tetragonal cell, II-30-7
 Thermal conductivity, II-2-8, II-12-2 of a gas, I-43-9 f
 Thermal equilibrium, I-41-3 ff
 Thermal ionization, I-42-5 ff
 Thermodynamics, I-39-2, I-45-1 ff, II-37-4 f laws of, I-44-1 ff
 Thompson, II-5-3
 Thompson atomic model, II-5-3
 Thompson scattering cross section, I-32-8
 Three-body problem, I-10-1
 Three-dimensional waves, II-20-8 f
 Thunderstorms, II-9-5 ff
 Tides, I-7-4 f
 Time, I-2-3, I-5-1 ff, I-8-1, I-8-2 retarded, I-28-2 standard of, I-5-5 transformation of, I-15-5 ff
 Torque, I-18-4, I-20-1 ff
 Torsion bar, II-38-5 ff
 Total internal reflection, II-33-12 f
 Transformation, Fourier, I-25-4 Galilean, I-12-11 linear, I-11-6 Lorentz, I-15-3, I-17-1, I-34-8, I-52-2, II-25-1, II-26-1 ff of time, I-15-5 ff of velocity, I-16-4 ff
 Transformer, II-16-4 f
 Transient, I-24-1 ff electrical, I-24-5 f
 Transient response, I-21-6
 Translation of axes, I-11-1 ff
 Transmission line, II-24-1 ff
 Transmitted waves, II-33-7 ff
 Travelling field, II-18-5 ff
 Triclinic lattice, II-30-7
 Trigonal lattice, II-30-7
 Twin paradox, I-16-3 f

- Two-dimensional field, II-7-2 ff
- Tycho Brahe, I-7-1
- Ultraviolet radiation, I-26-1
- Uncertainty principle, I-2-6, I-6-10 f, I-37-9, I-37-11, I-38-8 f, II-5-3
- Unit cell, I-38-5
- Unit vector, I-11-10, II-2-3
- Unworldliness, II-25-10
- van de Graaff generator, II-5-9, II-8-7
- Vector, I-11-5 ff
- Vector algebra, I-11-6 f
- Vector analysis, I-11-5, I-52-2
- Vector field, II-1-4 f, II-2-1 ff
 - flux of, II-3-2 ff
- Vector integrals, II-3-1 f
- Vector operator, II-2-6
- Vector potential, II-4-1 ff, II-15-1 ff
- Vector product, I-20-4
- Velocity, I-8-3, I-9-2 f
 - components of, I-9-3
 - transformation of, I-16-4 ff
- Velocity potential, II-12-9
- Vinci, Leonardo da, I-36-2
- Virtual work, principle of, I-4-5
- Viscosity, II-41-1 ff
 - coefficient of, II-41-2
- Viscous flow, II-41-4 f
- Vision, I-36-1 ff
 - binocular, I-36-4
 - color, I-35-1 ff
- Visual cortex, I-36-4
- Visual purple, I-35-9
- Voltmeter, II-16-1
- Volume strain, II-38-3
- Volume stress, II-38-3
- von Neumann, J., II-40-3
- Vortex lines, II-40-10 ff
- Vorticity, II-40-5
- Wall energy, II-37-6
- Wapstra, I-52-10
- Watt (unit), I-13-3
- Wave, I-51-1 ff, II-20-1 ff
 - electromagnetic, II-21-1 f
 - light, I-48-1
 - plane, II-20-1 ff
 - reflected, II-33-7 ff
 - shear, I-51-4, II-38-8
 - sinusoidal, I-29-2 f
 - spherical, II-20-12 ff, II-21-2 ff
 - three-dimensional, II-20-8 f
 - transmitted, II-33-7 ff
- Wave equation, I-47-1 ff, II-18-9 ff
- Wavefront, I-47-3
- Waveguides, II-24-1 ff
- Wavelength, I-19-3, I-26-1
- Wave number, I-29-2
- Weber, II-16-2
- Weber (unit), II-13-1
- "Wet" water, II-41-1 ff
- Weyl, H., I-11-1
- Wheeler, II-28-8
- Wilson, C. T. R., II-9-9
- Work, I-13-1 ff, I-14-1 ff
- X-rays, I-2-5, I-26-1
- X-ray diffraction, II-30-1
- Young, I-35-7
- Young's modulus, II-38-2
- Yukawa, H., I-2-8, II-28-13
- Yukawa potential, II-28-13
- Yustova, I-35-8
- Zeno, I-8-3
- Zero, absolute, I-1-5
- Zero curl, II-3-10 f, II-4-1
- Zero divergence, II-3-10 f, II-4-1
- Zero mass, I-2-10

Feynman's Preface



These are the lectures in physics that I gave last year and the year before to the freshman and sophomore classes at Caltech. The lectures are, of course, not verbatim—they have been edited, sometimes extensively and sometimes less so. The lectures form only part of the complete course. The whole group of 180 students gathered in a big lecture room twice a week to hear these lectures and then they broke up into small groups of 15 to 20 students in recitation sections under the guidance of a teaching assistant. In addition, there was a laboratory session once a week.

The special problem we tried to get at with these lectures was to maintain the interest of the very enthusiastic and rather smart students coming out of the high schools and into Caltech. They have heard a lot about how interesting and exciting physics is—the theory of relativity, quantum mechanics, and other modern ideas. By the end of two years of our previous course, many would be very discouraged because there were really very few grand, new, modern ideas presented to them. They were made to study inclined planes, electrostatics, and so forth, and after two years it was quite stultifying. The problem was whether or not we could make a course which would save the more advanced and excited student by maintaining his enthusiasm.

The lectures here are not in any way meant to be a survey course, but are very serious. I thought to address them to the most intelligent in the class and to make sure, if possible, that even the most intelligent student was unable to completely encompass everything that was in the lectures—by putting in suggestions of applications of the ideas and concepts in various directions outside the main line of attack. For this reason, though, I tried very hard to make all the statements as accurate as possible, to point out in every case where the equations and ideas fitted into the body of physics, and how—when they learned more—things would be modified. I also felt that for such students it is important to indicate what it is that they should—if they are sufficiently clever—be able to understand by deduction from what has been said before, and what is being put in as something new. When new ideas came in, I would try either to deduce them if they were deducible, or to explain that it *was* a new idea which hadn't any basis in terms of things they had already learned and which was not supposed to be provable—but was just added in.

At the start of these lectures, I assumed that the students knew something when they came out of high school—such things as geometrical optics, simple chemistry ideas, and so on. I also didn't see that there was any reason to make the lectures

Foreword

A great triumph of twentieth-century physics, the theory of quantum mechanics, is now nearly 40 years old, yet we have generally been giving our students their introductory course in physics (for many students, their last) with hardly more than a casual allusion to this central part of our knowledge of the physical world. We should do better by them. These lectures are an attempt to present them with the basic and essential ideas of the quantum mechanics in a way that would, hopefully, be comprehensible. The approach you will find here is novel, particularly at the level of a sophomore course, and was considered very much an experiment. After seeing how easily some of the students take to it, however, I believe that the experiment was a success. There is, of course, room for improvement, and it will come with more experience in the classroom. What you will find here is a record of that first experiment.

In the two-year sequence of the Feynman Lectures on Physics which were given from September 1961 through May 1963 for the introductory physics course at Caltech, the concepts of quantum physics were brought in whenever they were necessary for an understanding of the phenomena being described. In addition, the last twelve lectures of the second year were given over to a more coherent introduction to some of the concepts of quantum mechanics. It became clear as the lectures drew to a close, however, that not enough time had been left for the quantum mechanics. As the material was prepared, it was continually discovered that other important and interesting topics could be treated with the elementary tools that had been developed. There was also a fear that the too brief treatment of the Schrödinger wave function which had been included in the twelfth lecture would not provide a sufficient bridge to the more conventional treatments of many books the students might hope to read. It was therefore decided to extend the series with seven additional lectures; they were given to the sophomore class in May of 1964. These lectures rounded out and extended somewhat the material developed in the earlier lectures.

In this volume we have put together the lectures from both years with some adjustment of the sequence. In addition, two lectures originally given to the freshman class as an introduction to quantum physics have been lifted bodily from Volume I (where they were Chapters 37 and 38) and placed as the first two chapters here—to make this volume a self-contained unit, relatively independent of the first two. A few ideas about the quantization of angular momentum (including a discussion of the Stern-Gerlach experiment) had been introduced in Chapters 34 and 35 of Volume II, and familiarity with them is assumed; for the convenience of those who will not have that volume at hand, those two chapters are reproduced here as an Appendix.

This set of lectures tries to elucidate from the beginning those features of the quantum mechanics which are most basic and most general. The first lectures tackle head on the ideas of a probability amplitude, the interference of amplitudes, the abstract notion of a state, and the superposition and resolution of states—and the Dirac notation is used from the start. In each instance the ideas are introduced together with a detailed discussion of some specific examples—to try to make the physical ideas as real as possible. The time dependence of states including states of definite energy comes next, and the ideas are applied at once to the study of two-state systems. A detailed discussion of the ammonia maser provides the frame-

work for the introduction to radiation absorption and induced transitions. The lectures then go on to consider more complex systems, leading to a discussion of the propagation of electrons in a crystal, and to a rather complete treatment of the quantum mechanics of angular momentum. Our introduction to quantum mechanics ends in Chapter 20 with a discussion of the Schrödinger wave function, its differential equation, and the solution for the hydrogen atom.

The last chapter of this volume is not intended to be a part of the “course.” It is a “seminar” on superconductivity and was given in the spirit of some of the entertainment lectures of the first two volumes, with the intent of opening to the students a broader view of the relation of what they were learning to the general culture of physics. Feynman’s “epilogue” serves as the period to the three-volume series.

As explained in the Foreword to Volume I, these lectures were but one aspect of a program for the development of a new introductory course carried out at the California Institute of Technology under the supervision of the Physics Course Revision Committee (Robert Leighton, Victor Neher, and Matthew Sands). The program was made possible by a grant from the Ford Foundation. Many people helped with the technical details of the preparation of this volume: Marylou Clayton, Julie Curcio, James Hartle, Tom Harvey, Martin Israel, Patricia Preuss, Fanny Warren, and Barbara Zimmerman. Professors Gerry Neugebauer and Charles Wilts contributed greatly to the accuracy and clarity of the material by reviewing carefully much of the manuscript.

But the story of quantum mechanics you will find here is Richard Feynman’s. Our labors will have been well spent if we have been able to bring to others even some of the intellectual excitement we experienced as we saw the ideas unfold in his real-life Lectures on Physics.

December, 1964

MATTHEW SANDS

Contents

CHAPTER 1. QUANTUM BEHAVIOR

- 1-1 Atomic mechanics 1-1
- 1-2 An experiment with bullets 1-1
- 1-3 An experiment with waves 1-3
- 1-4 An experiment with electrons 1-4
- 1-5 The interference of electron waves 1-5
- 1-6 Watching the electrons 1-6
- 1-7 First principles of quantum mechanics 1-9
- 1-8 The uncertainty principle 1-11

CHAPTER 2. THE RELATION OF WAVE AND PARTICLE VIEWPOINTS

- 2-1 Probability wave amplitudes 2-1
- 2-2 Measurement of position and momentum 2-2
- 2-3 Crystal diffraction 2-4
- 2-4 The size of an atom 2-5
- 2-5 Energy levels 2-7
- 2-6 Philosophical implications 2-8

CHAPTER 3. PROBABILITY AMPLITUDES

- 3-1 The laws of combining amplitudes 3-1
- 3-2 The two-slit interference pattern 3-5
- 3-3 Scattering from a crystal 3-7
- 3-4 Identical particles 3-9

CHAPTER 4. IDENTICAL PARTICLES

- 4-1 Bose particles and Fermi particles 4-1
- 4-2 States with two Bose particles 4-3
- 4-3 States with n Bose particles 4-6
- 4-4 Emission and absorption of photons 4-7
- 4-5 The blackbody spectrum 4-8
- 4-6 Liquid helium 4-12
- 4-7 The exclusion principle 4-12

CHAPTER 5. SPIN ONE

- 5-1 Filtering atoms with a Stern-Gerlach apparatus 5-1
- 5-2 Experiments with filtered atoms 5-5
- 5-3 Stern-Gerlach filters in series 5-6
- 5-4 Base states 5-8
- 5-5 Interfering amplitudes 5-10
- 5-6 The machinery of quantum mechanics 5-12
- 5-7 Transforming to a different base 5-15
- 5-8 Other situations 5-16

CHAPTER 6. SPIN ONE-HALF

- 6-1 Transforming amplitudes 6-1
- 6-2 Transforming to a rotated coordinate system 6-3
- 6-3 Rotations about the z -axis 6-6
- 6-4 Rotations of 180° and 90° about y 6-9
- 6-5 Rotations about x 6-11
- 6-6 Arbitrary rotations 6-12

CHAPTER 7. THE DEPENDENCE OF AMPLITUDES ON TIME

- 7-1 Atoms at rest; stationary states 7-1
- 7-2 Uniform motion 7-3
- 7-3 Potential energy; energy conservation 7-6
- 7-4 Forces; the classical limit 7-9
- 7-5 The "precession" of a spin one-half particle 7-10

CHAPTER 8. THE HAMILTONIAN MATRIX

- 8-1 Amplitudes and vectors 8-1
- 8-2 Resolving state vectors 8-3
- 8-3 What are the base states of the world? 8-5
- 8-4 How states change with time 8-7
- 8-5 The Hamiltonian matrix 8-10
- 8-6 The ammonia molecule 8-11

CHAPTER 9. THE AMMONIA MASER

- 9-1 The states of an ammonia molecule 9-1
- 9-2 The molecule in a static electric field 9-5
- 9-3 Transitions in a time-dependent field 9-9
- 9-4 Transitions at resonance 9-11
- 9-5 Transitions off resonance 9-13
- 9-6 The absorption of light 9-14

CHAPTER 10. OTHER TWO-STATE SYSTEMS

- 10-1 The hydrogen molecular ion 10-1
- 10-2 Nuclear forces 10-6
- 10-3 The hydrogen molecule 10-8
- 10-4 The benzene molecule 10-10
- 10-5 Dyes 10-12
- 10-6 The Hamiltonian of a spin one-half particle in a magnetic field 10-12
- 10-7 The spinning electron in a magnetic field 10-15

CHAPTER 11. MORE TWO-STATE SYSTEMS

- 11-1 The Pauli spin matrices 11-1
- 11-2 The spin matrices as operators 11-5
- 11-3 The solution of the two-state equations 11-8
- 11-4 The polarization states of the photon 11-9
- 11-5 The neutral K-meson 11-12
- 11-6 Generalization to N -state systems 11-20

CHAPTER 12. THE HYPERFINE SPLITTING IN HYDROGEN

- 12-1 Base states for a system with two spin one-half particles 12-1
- 12-2 The Hamiltonian for the ground state of hydrogen 12-3
- 12-3 The energy levels 12-7
- 12-4 The Zeeman splitting 12-9
- 12-5 The states in a magnetic field 12-12
- 12-6 The projection matrix for spin one 12-14

CHAPTER 13. PROPAGATION IN A CRYSTAL LATTICE

- 13-1 States for an electron in a one-dimensional lattice 13-1
- 13-2 States of definite energy 13-3
- 13-3 Time-dependent states 13-6
- 13-4 An electron in a three-dimensional lattice 13-7
- 13-5 Other states in a lattice 13-8
- 13-6 Scattering by imperfections in the lattice 13-10
- 13-7 Trapping by a lattice imperfection 13-12
- 13-8 Scattering amplitudes and bound states 13-13

CHAPTER 14. SEMICONDUCTORS

- 14-1 Electrons and holes in semiconductors 14-1
- 14-2 Impure semiconductors 14-4
- 14-3 The Hall effect 14-7
- 14-4 Semiconductor junctions 14-8
- 14-5 Rectification at a semiconductor junction 14-10
- 14-6 The transistor 14-11

CHAPTER 15. THE INDEPENDENT PARTICLE APPROXIMATION

- 15-1 Spin waves 15-1
- 15-2 Two spin waves 15-4
- 15-3 Independent particles 15-6
- 15-4 The benzene molecule 15-7
- 15-5 More organic chemistry 15-10
- 15-6 Other uses of the approximation 15-12

CHAPTER 16. THE DEPENDENCE OF AMPLITUDES ON POSITION

- 16-1 Amplitudes on a line 16-1
- 16-2 The wave function 16-5
- 16-3 States of definite momentum 16-7
- 16-4 Normalization of the states in x 16-9
- 16-5 The Schrödinger equation 16-11
- 16-6 Quantized energy levels 16-14

CHAPTER 17. SYMMETRY AND CONSERVATION LAWS

- 17-1 Symmetry 17-1
- 17-2 Symmetry and conservation 17-3
- 17-3 The conservation laws 17-7
- 17-4 Polarized light 17-9
- 17-5 The disintegration of the Λ^0 17-11
- 17-6 Summary of the rotation matrices 17-15

CHAPTER 18. ANGULAR MOMENTUM

- 18-1 Electric dipole radiation 18-1
- 18-2 Light scattering 18-3
- 18-3 The annihilation of positronium 18-5
- 18-4 Rotation matrix for any spin 18-9
- 18-5 Measuring a nuclear spin 18-13
- 18-6 Composition of angular momentum 18-14
- Added Note 1: Derivation of the rotation matrix 18-19
- Added Note 2: Conservation of parity in photon emission 18-22

CHAPTER 19. THE HYDROGEN ATOM AND THE PERIODIC TABLE

- 19-1 Schrödinger's equation for the hydrogen atom 19-1
- 19-2 Spherically symmetric solutions 19-2
- 19-3 States with an angular dependence 19-6
- 19-4 The general solution for hydrogen 19-10
- 19-5 The hydrogen wave functions 19-12
- 19-6 The periodic table 19-13

CHAPTER 20. OPERATORS

- 20-1 Operations and operators 20-1
- 20-2 Average energies 20-3
- 20-3 The average energy of an atom 20-6
- 20-4 The position operator 20-8
- 20-5 The momentum operator 20-9
- 20-6 Angular momentum 20-14
- 20-7 The change of averages with time 20-15

CHAPTER 21. THE SCHRÖDINGER EQUATION IN A CLASSICAL CONTEXT: A SEMINAR ON SUPERCONDUCTIVITY

- 21-1 Schrödinger's equation in a magnetic field 21-1
- 21-2 The equation of continuity for probabilities 21-3
- 21-3 Two kinds of momentum 21-4
- 21-4 The meaning of the wave function 21-6
- 21-5 Superconductivity 21-7
- 21-6 The Meissner effect 21-8
- 21-7 Flux quantization 21-10
- 21-8 The dynamics of superconductivity 21-12
- 21-9 The Josephson junction 21-14

FEYNMAN'S EPILOGUE

APPENDIX

INDEX

Quantum Behavior

1-1 Atomic mechanics

“Quantum mechanics” is the description of the behavior of matter and light in all its details and, in particular, of the happenings on an atomic scale. Things on a very small scale behave like nothing that you have any direct experience about. They do not behave like waves, they do not behave like particles, they do not behave like clouds, or billiard balls, or weights on springs, or like anything that you have ever seen.

Newton thought that light was made up of particles, but then it was discovered that it behaves like a wave. Later, however (in the beginning of the twentieth century), it was found that light did indeed sometimes behave like a particle. Historically, the electron, for example, was thought to behave like a particle, and then it was found that in many respects it behaved like a wave. So it really behaves like neither. Now we have given up. We say: “It is like *neither*.”

There is one lucky break, however—electrons behave just like light. The quantum behavior of atomic objects (electrons, protons, neutrons, photons, and so on) is the same for all, they are all “particle waves,” or whatever you want to call them. So what we learn about the properties of electrons (which we shall use for our examples) will apply also to all “particles,” including photons of light.

The gradual accumulation of information about atomic and small-scale behavior during the first quarter of this century, which gave some indications about how small things do behave, produced an increasing confusion which was finally resolved in 1926 and 1927 by Schrödinger, Heisenberg, and Born. They finally obtained a consistent description of the behavior of matter on a small scale. We take up the main features of that description in this chapter.

Because atomic behavior is so unlike ordinary experience, it is very difficult to get used to, and it appears peculiar and mysterious to everyone—both to the novice and to the experienced physicist. Even the experts do not understand it the way they would like to, and it is perfectly reasonable that they should not, because all of direct, human experience and of human intuition applies to large objects. We know how large objects will act, but things on a small scale just do not act that way. So we have to learn about them in a sort of abstract or imaginative fashion and not by connection with our direct experience.

In this chapter we shall tackle immediately the basic element of the mysterious behavior in its most strange form. We choose to examine a phenomenon which is impossible, *absolutely* impossible, to explain in any classical way, and which has in it the heart of quantum mechanics. In reality, it contains the *only* mystery. We cannot make the mystery go away by “explaining” how it works. We will just *tell* you how it works. In telling you how it works we will have told you about the basic peculiarities of all quantum mechanics.

1-2 An experiment with bullets

To try to understand the quantum behavior of electrons, we shall compare and contrast their behavior, in a particular experimental setup, with the more familiar behavior of particles like bullets, and with the behavior of waves like water waves. We consider first the behavior of bullets in the experimental setup shown diagrammatically in Fig. 1-1. We have a machine gun that shoots a stream of bullets. It is not a very good gun, in that it sprays the bullets (randomly) over a fairly large angular spread, as indicated in the figure. In front of the gun we have

1-1 Atomic mechanics

1-2 An experiment with bullets

1-3 An experiment with waves

1-4 An experiment with electrons

1-5 The interference of electron waves

1-6 Watching the electrons

1-7 First principles of quantum mechanics

1-8 The uncertainty principle

Note: This chapter is almost exactly the same as Chapter 37 of Volume I.

a wall (made of armor plate) that has in it two holes just about big enough to let a bullet through. Beyond the wall is a backstop (say a thick wall of wood) which will “absorb” the bullets when they hit it. In front of the wall we have an object which we shall call a “detector” of bullets. It might be a box containing sand. Any bullet that enters the detector will be stopped and accumulated. When we wish, we can empty the box and count the number of bullets that have been caught. The detector can be moved back and forth (in what we will call the x -direction). With this apparatus, we can find out experimentally the answer to the question: “What is the probability that a bullet which passes through the holes in the wall will arrive at the backstop at the distance x from the center?” First, you should realize that we should talk about probability, because we cannot say definitely where any particular bullet will go. A bullet which happens to hit one of the holes may bounce off the edges of the hole, and may end up anywhere at all. By “probability” we mean the chance that the bullet will arrive at the detector, which we can measure by counting the number which arrive at the detector in a certain time and then taking the ratio of this number to the *total* number that hit the backstop during that time. Or, if we assume that the gun always shoots at the same rate during the measurements, the probability we want is just proportional to the number that reach the detector in some standard time interval.

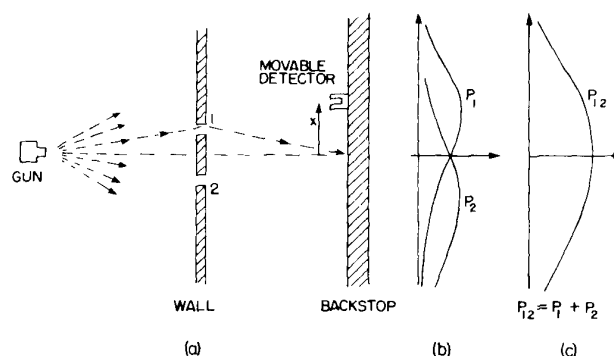


Fig. 1-1. Interference experiment with bullets.

For our present purposes we would like to imagine a somewhat idealized experiment in which the bullets are not real bullets, but are *indestructible* bullets—they cannot break in half. In our experiment we find that bullets always arrive in lumps, and when we find something in the detector, it is always one whole bullet. If the rate at which the machine gun fires is made very low, we find that at any given moment either nothing arrives, or one and only one—exactly one—bullet arrives at the backstop. Also, the size of the lump certainly does not depend on the rate of firing of the gun. We shall say: “Bullets *always* arrive in identical lumps.” What we measure with our detector is the probability of arrival of a lump. And we measure the probability as a function of x . The result of such measurements with this apparatus (we have not yet done the experiment, so we are really imagining the result) are plotted in the graph drawn in part (c) of Fig. 1-1. In the graph we plot the probability to the right and x vertically, so that the x -scale fits the diagram of the apparatus. We call the probability P_{12} because the bullets may have come either through hole 1 or through hole 2. You will not be surprised that P_{12} is large near the middle of the graph but gets small if x is very large. You may wonder, however, why P_{12} has its maximum value at $x = 0$. We can understand this fact if we do our experiment again after covering up hole 2, and once more while covering up hole 1. When hole 2 is covered, bullets can pass only through hole 1, and we get the curve marked P_1 in part (b) of the figure. As you would expect, the maximum of P_1 occurs at the value of x which is on a straight line with the gun and hole 1. When hole 1 is closed, we get the symmetric curve P_2 drawn in the figure. P_2 is the probability distribution for bullets that pass through hole 2. Comparing parts (b) and (c) of Fig. 1-1, we find the important result that

$$P_{12} = P_1 + P_2. \quad (1.1)$$

The probabilities just add together. The effect with both holes open is the sum of the effects with each hole open alone. We shall call this result an observation of “no interference,” for a reason that you will see later. So much for bullets. They come in lumps, and their probability of arrival shows no interference.

1-3 An experiment with waves

Now we wish to consider an experiment with water waves. The apparatus is shown diagrammatically in Fig. 1-2. We have a shallow trough of water. A small object labeled the “wave source” is jigged up and down by a motor and makes circular waves. To the right of the source we have again a wall with two holes, and beyond that is a second wall, which, to keep things simple, is an “absorber,” so that there is no reflection of the waves that arrive there. This can be done by building a gradual sand “beach.” In front of the beach we place a detector which can be moved back and forth in the x -direction, as before. The detector is now a device which measures the “intensity” of the wave motion. You can imagine a gadget which measures the height of the wave motion, but whose scale is calibrated in proportion to the *square* of the actual height, so that the reading is proportional to the intensity of the wave. Our detector reads, then, in proportion to the *energy* being carried by the wave—or rather, the rate at which energy is carried to the detector.

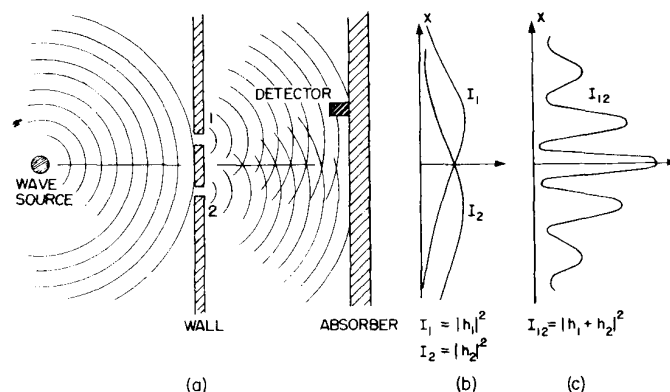


Fig. 1-2. Interference experiment with water waves.

With our wave apparatus, the first thing to notice is that the intensity can have *any* size. If the source just moves a very small amount, then there is just a little bit of wave motion at the detector. When there is more motion at the source, there is more intensity at the detector. The intensity of the wave can have any value at all. We would *not* say that there was any “lumpiness” in the wave intensity.

Now let us measure the wave intensity for various values of x (keeping the wave source operating always in the same way). We get the interesting-looking curve marked I_{12} in part (c) of the figure.

We have already worked out how such patterns can come about when we studied the interference of electric waves in Volume I. In this case we would observe that the original wave is diffracted at the holes, and new circular waves spread out from each hole. If we cover one hole at a time and measure the intensity distribution at the absorber we find the rather simple intensity curves shown in part (b) of the figure. I_1 is the intensity of the wave from hole 1 (which we find by measuring when hole 2 is blocked off) and I_2 is the intensity of the wave from hole 2 (seen when hole 1 is blocked).

The intensity I_{12} observed when both holes are open is certainly *not* the sum of I_1 and I_2 . We say that there is “interference” of the two waves. At some places (where the curve I_{12} has its maxima) the waves are “in phase” and the wave peaks add together to give a large amplitude and, therefore, a large intensity. We say that the two waves are “interfering constructively” at such places. There will be such constructive interference wherever the distance from the detector to one hole is a whole number of wavelengths larger (or shorter) than the distance from the detector to the other hole.

At those places where the two waves arrive at the detector with a phase difference of π (where they are “out of phase”) the resulting wave motion at the detector will be the difference of the two amplitudes. The waves “interfere destructively,” and we get a low value for the wave intensity. We expect such low values wherever the distance between hole 1 and the detector is different from the distance between hole 2 and the detector by an odd number of half-wavelengths. The low values of I_{12} in Fig. 1-2 correspond to the places where the two waves interfere destructively.

You will remember that the quantitative relationship between I_1 , I_2 , and I_{12} can be expressed in the following way: The instantaneous height of the water wave at the detector for the wave from hole 1 can be written as (the real part of) $h_1 e^{i\omega t}$, where the “amplitude” h_1 is, in general, a complex number. The intensity is proportional to the mean squared height or, when we use the complex numbers, to the absolute value squared $|h_1|^2$. Similarly, for hole 2 the height is $h_2 e^{i\omega t}$ and the intensity is proportional to $|h_2|^2$. When both holes are open, the wave heights add to give the height $(h_1 + h_2)e^{i\omega t}$ and the intensity $|h_1 + h_2|^2$. Omitting the constant of proportionality for our present purposes, the proper relations for *interfering waves* are

$$I_1 = |h_1|^2, \quad I_2 = |h_2|^2, \quad I_{12} = |h_1 + h_2|^2. \quad (1.2)$$

You will notice that the result is quite different from that obtained with bullets (Eq. 1-1). If we expand $|h_1 + h_2|^2$ we see that

$$|h_1 + h_2|^2 = |h_1|^2 + |h_2|^2 + 2|h_1||h_2|\cos\delta, \quad (1.3)$$

where δ is the phase difference between h_1 and h_2 . In terms of the intensities, we could write

$$I_{12} = I_1 + I_2 + 2\sqrt{I_1 I_2} \cos\delta. \quad (1.4)$$

The last term in (1.4) is the “interference term.” So much for water waves. The intensity can have any value, and it shows interference.

1-4 An experiment with electrons

Now we imagine a similar experiment with electrons. It is shown diagrammatically in Fig. 1-3. We make an electron gun which consists of a tungsten wire heated by an electric current and surrounded by a metal box with a hole in it. If the wire is at a negative voltage with respect to the box, electrons emitted by the wire will be accelerated toward the walls and some will pass through the hole. All the electrons which come out of the gun will have (nearly) the same energy. In front of the gun is again a wall (just a thin metal plate) with two holes in it. Beyond the wall is another plate which will serve as a “backstop.” In front of the backstop we place a movable detector. The detector might be a geiger counter or, perhaps better, an electron multiplier, which is connected to a loudspeaker.

We should say right away that you should not try to set up this experiment (as you could have done with the two we have already described). This experiment

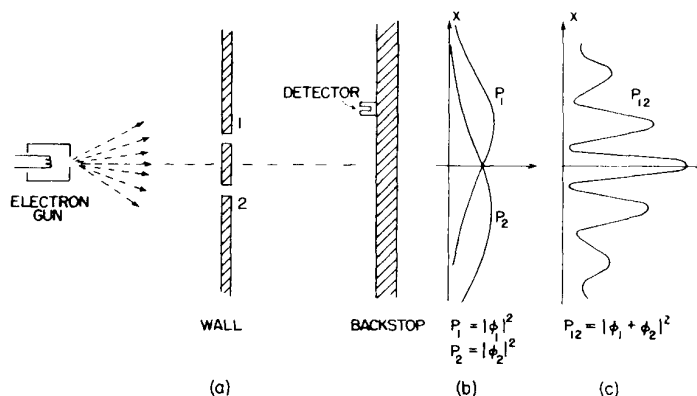


Fig. 1-3. Interference experiment with electrons.

has never been done in just this way. The trouble is that the apparatus would have to be made on an impossibly small scale to show the effects we are interested in. We are doing a “thought experiment,” which we have chosen because it is easy to think about. We know the results that *would* be obtained because there *are* many experiments that have been done, in which the scale and the proportions have been chosen to show the effects we shall describe.

The first thing we notice with our electron experiment is that we hear sharp “clicks” from the detector (that is, from the loudspeaker). And all “clicks” are the same. There are *no* “half-clicks.”

We would also notice that the “clicks” come very erratically. Something like: click click-click . . . click click click-click click . . . , etc., just as you have, no doubt, heard a geiger counter operating. If we count the clicks which arrive in a sufficiently long time—say for many minutes—and then count again for another equal period, we find that the two numbers are very nearly the same. So we can speak of the *average rate* at which the clicks are heard (so-and-so-many clicks per minute on the average).

As we move the detector around, the *rate* at which the clicks appear is faster or slower, but the size (loudness) of each click is always the same. If we lower the temperature of the wire in the gun, the rate of clicking slows down, but still each click sounds the same. We would notice also that if we put two separate detectors at the backstop, one *or* the other would click, but never both at once. (Except that once in a while, if there were two clicks very close together in time, our ear might not sense the separation.) We conclude, therefore, that whatever arrives at the backstop arrives in “lumps.” All the “lumps” are the same size: only whole “lumps” arrive, and they arrive one at a time at the backstop. We shall say: “Electrons always arrive in identical lumps.”

Just as for our experiment with bullets, we can now proceed to find experimentally the answer to the question: “What is the relative probability that an electron ‘lump’ will arrive at the backstop at various distances x from the center?” As before, we obtain the relative probability by observing the rate of clicks, holding the operation of the gun constant. The probability that lumps will arrive at a particular x is proportional to the average rate of clicks at that x .

The result of our experiment is the interesting curve marked P_{12} in part (c) of Fig. 1-3. Yes! That is the way electrons go.

1-5 The interference of electron waves

Now let us try to analyze the curve of Fig. 1-3 to see whether we can understand the behavior of the electrons. The first thing we would say is that since they come in lumps, each lump, which we may as well call an electron, has come either through hole 1 or through hole 2. Let us write this in the form of a “Proposition”:

Proposition A: Each electron *either* goes through hole 1 *or* it goes through hole 2.

Assuming Proposition A, all electrons that arrive at the backstop can be divided into two classes: (1) those that come through hole 1, and (2) those that come through hole 2. So our observed curve must be the sum of the effects of the electrons which come through hole 1 and the electrons which come through hole 2. Let us check this idea by experiment. First, we will make a measurement for those electrons that come through hole 1. We block off hole 2 and make our counts of the clicks from the detector. From the clicking rate, we get P_1 . The result of the measurement is shown by the curve marked P_1 in part (b) of Fig. 1-3. The result seems quite reasonable. In a similar way, we measure P_2 , the probability distribution for the electrons that come through hole 2. The result of this measurement is also drawn in the figure.

The result P_{12} obtained with *both* holes open is clearly not the sum of P_1 and P_2 , the probabilities for each hole alone. In analogy with our water-wave experi-

ment, we say: "There is interference."

$$\text{For electrons: } P_{12} \neq P_1 + P_2. \quad (1.5)$$

How can such an interference come about? Perhaps we should say: "Well, that means, presumably, that it is *not true* that the lumps go either through hole 1 or hole 2, because if they did, the probabilities should add. Perhaps they go in a more complicated way. They split in half and . . ." But no! They cannot, they always arrive in lumps . . . "Well, perhaps some of them go through 1, and then they go around through 2, and then around a few more times, or by some other complicated path . . . then by closing hole 2, we changed the chance that an electron that *started out* through hole 1 would finally get to the backstop . . ." But notice! There are some points at which very few electrons arrive when *both* holes are open, but which receive many electrons if we close one hole, so *closing* one hole *increased* the number from the other. Notice, however, that at the center of the pattern, P_{12} is more than twice as large as $P_1 + P_2$. It is as though closing one hole *decreased* the number of electrons which come through the other hole. It seems hard to explain *both* effects by proposing that the electrons travel in complicated paths.

It is all quite mysterious. And the more you look at it the more mysterious it seems. Many ideas have been concocted to try to explain the curve for P_{12} in terms of individual electrons going around in complicated ways through the holes. None of them has succeeded. None of them can get the right curve for P_{12} in terms of P_1 and P_2 .

Yet, surprisingly enough, the *mathematics* for relating P_1 and P_2 to P_{12} is extremely simple. For P_{12} is just like the curve I_{12} of Fig. 1-2, and *that* was simple. What is going on at the backstop can be described by two complex numbers that we can call ϕ_1 and ϕ_2 (they are functions of x , of course). The absolute square of ϕ_1 gives the effect with only hole 1 open. That is, $P_1 = |\phi_1|^2$. The effect with only hole 2 open is given by ϕ_2 in the same way. That is, $P_2 = |\phi_2|^2$. And the combined effect of the two holes is just $P_{12} = |\phi_1 + \phi_2|^2$. The *mathematics* is the same as that we had for the water waves! (It is hard to see how one could get such a simple result from a complicated game of electrons going back and forth through the plate on some strange trajectory.)

We conclude the following: The electrons arrive in lumps, like particles, and the probability of arrival of these lumps is distributed like the distribution of intensity of a wave. It is in this sense that an electron behaves "sometimes like a particle and sometimes like a wave."

Incidentally, when we were dealing with classical waves we defined the intensity as the mean over time of the square of the wave amplitude, and we used complex numbers as a mathematical trick to simplify the analysis. But in quantum mechanics it turns out that the amplitudes *must* be represented by complex numbers. The real parts alone will not do. That is a technical point, for the moment, because the formulas look just the same.

Since the probability of arrival through both holes is given so simply, although it is not equal to $(P_1 + P_2)$, that is really all there is to say. But there are a large number of subtleties involved in the fact that nature does work this way. We would like to illustrate some of these subtleties for you now. First, since the number that arrives at a particular point is *not* equal to the number that arrives through 1 plus the number that arrives through 2, as we would have concluded from Proposition A, undoubtedly we should conclude that *Proposition A is false*. It is *not* true that the electrons go *either* through hole 1 or hole 2. But that conclusion can be tested by another experiment.

1-6 Watching the electrons

We shall now try the following experiment. To our electron apparatus we add a very strong light source, placed behind the wall and between the two holes, as shown in Fig. 1-4. We know that electric charges scatter light. So when an

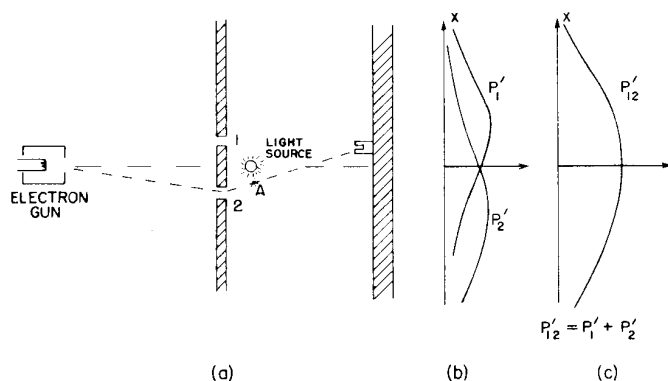


Fig. 1-4. A different electron experiment.

electron passes, however it does pass, on its way to the detector, it will scatter some light to our eye, and we can *see* where the electron goes. If, for instance, an electron were to take the path via hole 2 that is sketched in Fig. 1-4, we should see a flash of light coming from the vicinity of the place marked *A* in the figure. If an electron passes through hole 1, we would expect to see a flash from the vicinity of the upper hole. If it should happen that we get light from both places at the same time, because the electron divides in half . . . Let us just do the experiment!

Here is what we see: *every* time that we hear a “click” from our electron detector (at the backstop), we *also see* a flash of light *either* near hole 1 *or* near hole 2, but *never* both at once! And we observe the same result no matter where we put the detector. From this observation we conclude that when we look at the electrons we find that the electrons go either through one hole or the other. Experimentally, Proposition A is necessarily true.

What, then, is wrong with our argument *against* Proposition A? Why *isn't* P_{12} just equal to $P_1 + P_2$? Back to experiment! Let us keep track of the electrons and find out what they are doing. For each position (*x*-location) of the detector we will count the electrons that arrive and *also* keep track of which hole they went through, by watching for the flashes. We can keep track of things this way: whenever we hear a “click” we will put a count in Column 1 if we see the flash near hole 1, and if we see the flash near hole 2, we will record a count in Column 2. Every electron which arrives is recorded in one of two classes: those which come through 1 and those which come through 2. From the number recorded in Column 1 we get the probability P'_1 that an electron will arrive at the detector via hole 1; and from the number recorded in Column 2 we get P'_2 , the probability that an electron will arrive at the detector via hole 2. If we now repeat such a measurement for many values of *x*, we get the curves for P'_1 and P'_2 shown in part (b) of Fig. 1-4.

Well, that is not too surprising! We get for P'_1 something quite similar to what we got before for P_1 by blocking off hole 2; and P'_2 is similar to what we got by blocking hole 1. So there is *not* any complicated business like going through both holes. When we watch them, the electrons come through just as we would expect them to come through. Whether the holes are closed or open, those which we see come through hole 1 are distributed in the same way whether hole 2 is open or closed.

But wait! What do we have *now* for the *total* probability, the probability that an electron will arrive at the detector by any route? We already have that information. We just pretend that we never looked at the light flashes, and we lump together the detector clicks which we have separated into the two columns. We *must* just *add* the numbers. For the probability that an electron will arrive at the backstop by passing through *either* hole, we do find $P'_{12} = P'_1 + P'_2$. That is, although we succeeded in watching which hole our electrons come through, we no longer get the old interference curve P_{12} , but a new one, P'_{12} , showing no interference! If we turn out the light P_{12} is restored.

We must conclude that *when we look at the electrons* the distribution of them on the screen is different than when we do not look. Perhaps it is turning on our light source that disturbs things? It must be that the electrons are very delicate, and the light, when it scatters off the electrons, gives them a jolt that changes their

motion. We know that the electric field of the light acting on a charge will exert a force on it. So perhaps we *should* expect the motion to be changed. Anyway, the light exerts a big influence on the electrons. By trying to “watch” the electrons we have changed their motions. That is, the jolt given to the electron when the photon is scattered by it is such as to change the electron’s motion enough so that if it *might* have gone to where P_{12} was at a maximum it will instead land where P_{12} was a minimum; that is why we no longer see the wavy interference effects.

You may be thinking: “Don’t use such a bright source! Turn the brightness down! The light waves will then be weaker and will not disturb the electrons so much. Surely, by making the light dimmer and dimmer, eventually the wave will be weak enough that it will have a negligible effect.” O.K. Let’s try it. The first thing we observe is that the flashes of light scattered from the electrons as they pass by does *not* get weaker. *It is always the same-sized flash.* The only thing that happens as the light is made dimmer is that sometimes we hear a “click” from the detector but see *no flash at all*. The electron has gone by without being “seen.” What we are observing is that light *also* acts like electrons, we *knew* that it was “wavy,” but now we find that it is also “lumpy.” It always arrives—or is scattered—in lumps that we call “photons.” As we turn down the *intensity* of the light source we do not change the *size* of the photons, only the *rate* at which they are emitted. *That* explains why, when our source is dim, some electrons get by without being seen. There did not happen to be a photon around at the time the electron went through.

This is all a little discouraging. If it is true that whenever we “see” the electron we see the same-sized flash, then those electrons we see are *always* the disturbed ones. Let us try the experiment with a dim light anyway. Now whenever we hear a click in the detector we will keep a count in three columns: in Column (1) those electrons seen by hole 1, in Column (2) those electrons seen by hole 2, and in Column (3) those electrons not seen at all. When we work up our data (computing the probabilities) we find these results: Those “seen by hole 1” have a distribution like P'_1 ; those “seen by hole 2” have a distribution like P'_2 (so that those “seen by either hole 1 or 2” have a distribution like P'_{12}); and those “not seen at all” have a “wavy” distribution just like P_{12} of Fig. 1-3! *If the electrons are not seen, we have interference!*

That is understandable. When we do not see the electron, no photon disturbs it, and when we do see it, a photon has disturbed it. There is always the same amount of disturbance because the light photons all produce the same-sized effects and the effect of the photons being scattered is enough to smear out any interference effect.

Is there not *some* way we can see the electrons without disturbing them? We learned in an earlier chapter that the momentum carried by a “photon” is inversely proportional to its wavelength ($p = h/\lambda$). Certainly the jolt given to the electron when the photon is scattered toward our eye depends on the momentum that photon carries. Aha! If we want to disturb the electrons only slightly we should not have lowered the *intensity* of the light, we should have lowered its *frequency* (the same as increasing its wavelength). Let us use light of a redder color. We could even use infrared light, or radiowaves (like radar), and “see” where the electron went with the help of some equipment that can “see” light of these longer wavelengths. If we use “gentler” light perhaps we can avoid disturbing the electrons so much.

Let us try the experiment with longer waves. We shall keep repeating our experiment, each time with light of a longer wavelength. At first, nothing seems to change. The results are the same. Then a terrible thing happens. You remember that when we discussed the microscope we pointed out that, due to the *wave nature* of the light, there is a limitation on how close two spots can be and still be seen as two separate spots. This distance is of the order of the wavelength of light. So now, when we make the wavelength longer than the distance between our holes, we see a *big* fuzzy flash when the light is scattered by the electrons. We can no longer tell which hole the electron went through! We just know it went somewhere! And it is just with light of this color that we find that the jolts given to the electron

are small enough so that P'_{12} begins to look like P_{12} —that we begin to get some interference effect. And it is only for wavelengths much longer than the separation of the two holes (when we have no chance at all of telling where the electron went) that the disturbance due to the light gets sufficiently small that we again get the curve P_{12} shown in Fig. 1-3.

In our experiment we find that it is impossible to arrange the light in such a way that one can tell which hole the electron went through, and at the same time not disturb the pattern. It was suggested by Heisenberg that the then new laws of nature could only be consistent if there were some basic limitation on our experimental capabilities not previously recognized. He proposed, as a general principle, his *uncertainty principle*, which we can state in terms of our experiment as follows: "It is impossible to design an apparatus to determine which hole the electron passes through, that will not at the same time disturb the electrons enough to destroy the interference pattern." If an apparatus is capable of determining which hole the electron goes through, it *cannot* be so delicate that it does not disturb the pattern in an essential way. No one has ever found (or even thought of) a way around the uncertainty principle. So we must assume that it describes a basic characteristic of nature.

The complete theory of quantum mechanics which we now use to describe atoms and, in fact, all matter, depends on the correctness of the uncertainty principle. Since quantum mechanics is such a successful theory, our belief in the uncertainty principle is reinforced. But if a way to "beat" the uncertainty principle were ever discovered, quantum mechanics would give inconsistent results and would have to be discarded as a valid theory of nature.

"Well," you say, "what about Proposition A? Is it true, or is it *not* true, that the electron either goes through hole 1 or it goes through hole 2?" The only answer that can be given is that we have found from experiment that there is a certain special way that we have to think in order that we do not get into inconsistencies. What we must say (to avoid making wrong predictions) is the following. If one looks at the holes or, more accurately, if one has a piece of apparatus which is capable of determining whether the electrons go through hole 1 or hole 2, then one *can* say that it goes either through hole 1 or hole 2. *But*, when one does *not* try to tell which way the electron goes, when there is nothing in the experiment to disturb the electrons, then one may *not* say that an electron goes either through hole 1 or hole 2. If one does say that, and starts to make any deductions from the statement, he will make errors in the analysis. This is the logical tightrope on which we must walk if we wish to describe nature successfully.

If the motion of all matter—as well as electrons—must be described in terms of waves, what about the bullets in our first experiment? Why didn't we see an interference pattern there? It turns out that for the bullets the wavelengths were so tiny that the interference patterns became very fine. So fine, in fact, that with any detector of finite size one could not distinguish the separate maxima and minima. What we saw was only a kind of average, which is the classical curve. In Fig. 1-5 we have tried to indicate schematically what happens with large-scale objects. Part (a) of the figure shows the probability distribution one might predict for bullets, using quantum mechanics. The rapid wiggles are supposed to represent the interference pattern one gets for waves of very short wavelength. Any physical detector, however, straddles several wiggles of the probability curve, so that the measurements show the smooth curve drawn in part (b) of the figure.

1-7 First principles of quantum mechanics

We will now write a summary of the main conclusions of our experiments. We will, however, put the results in a form which makes them true for a general class of such experiments. We can write our summary more simply if we first define an "ideal experiment" as one in which there are no uncertain external influences, i.e., no jiggling or other things going on that we cannot take into ac-

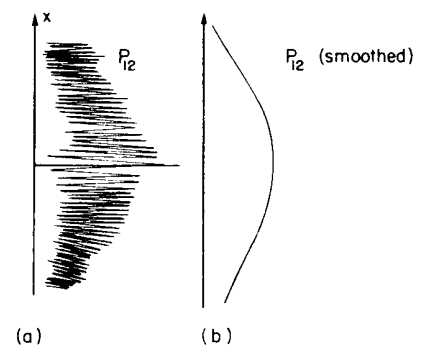


Fig. 1-5. Interference pattern with bullets: (a) actual (schematic), (b) observed.

count. We would be quite precise if we said: "An ideal experiment is one in which all of the initial and final conditions of the experiment are completely specified." What we will call "an event" is, in general, just a specific set of initial and final conditions. (For example: "an electron leaves the gun, arrives at the detector, and nothing else happens.") Now for our summary.

SUMMARY

- (1) The probability of an event in an ideal experiment is given by the square of the absolute value of a complex number ϕ which is called the probability amplitude:

$$\begin{aligned} P &= \text{probability,} \\ \phi &= \text{probability amplitude,} \\ P &= |\phi|^2 \end{aligned} \tag{1.6}$$

- (2) When an event can occur in several alternative ways, the probability amplitude for the event is the sum of the probability amplitudes for each way considered separately. There is interference:

$$\begin{aligned} \phi &= \phi_1 + \phi_2, \\ P &= |\phi_1 + \phi_2|^2 \end{aligned} \tag{1.7}$$

- (3) If an experiment is performed which is capable of determining whether one or another alternative is actually taken, the probability of the event is the sum of the probabilities for each alternative. The interference is lost:

$$P = P_1 + P_2. \tag{1.8}$$

One might still like to ask: "How does it work? What is the machinery behind the law?" No one has found any machinery behind the law. No one can "explain" any more than we have just "explained." No one will give you any deeper representation of the situation. We have no ideas about a more basic mechanism from which these results can be deduced.

We would like to emphasize a very important difference between classical and quantum mechanics. We have been talking about the probability that an electron will arrive in a given circumstance. We have implied that in our experimental arrangement (or even in the best possible one) it would be impossible to predict exactly what would happen. We can only predict the odds! This would mean, if it were true, that physics has given up on the problem of trying to predict exactly what will happen in a definite circumstance. Yes! physics *has* given up. *We do not know how to predict what would happen in a given circumstance*, and we believe now that it is impossible—that the only thing that can be predicted is the probability of different events. It must be recognized that this is a retrenchment in our earlier ideal of understanding nature. It may be a backward step, but no one has seen a way to avoid it.

We make now a few remarks on a suggestion that has sometimes been made to try to avoid the description we have given: "Perhaps the electron has some kind of internal works—some inner variables—that we do not yet know about. Perhaps that is why we cannot predict what will happen. If we could look more closely at the electron, we could be able to tell where it would end up." So far as we know, that is impossible. We would still be in difficulty. Suppose we were to assume that inside the electron there is some kind of machinery that determines where it is going to end up. That machine must *also* determine which hole it is going to go through on its way. But we must not forget that what is inside the electron should not be dependent on what *we* do, and in particular upon whether we open or close one of the holes. So if an electron, before it starts, has already made up its mind (a) which hole it is going to use, and (b) where it is going to land, we should find P_1 for those electrons that have chosen hole 1, P_2 for those that have chosen hole 2, and *necessarily* the sum $P_1 + P_2$ for those that arrive through the two holes. There seems to be no way around this. But we have verified experimentally that that is not the case. And no one has figured a way out of this puzzle. So at the

present time we must limit ourselves to computing probabilities. We say “at the present time,” but we suspect very strongly that it is something that will be with us forever—that it is impossible to beat that puzzle—that this is the way nature really *is*.

1-8 The uncertainty principle

This is the way Heisenberg stated the uncertainty principle originally: If you make the measurement on any object, and you can determine the x -component of its momentum with an uncertainty Δp , you cannot, at the same time, know its x -position more accurately than $\Delta x = h/\Delta p$, where h is a definite fixed number given by nature. It is called “Planck’s constant,” and is approximately 6.63×10^{-34} joule-seconds. The uncertainties in the position and momentum of a particle at any instant must have their product greater than Planck’s constant. This is a special case of the uncertainty principle that was stated above more generally. The more general statement was that one cannot design equipment in any way to determine which of two alternatives is taken, without, at the same time, destroying the pattern of interference.

Let us show for one particular case that the kind of relation given by Heisenberg must be true in order to keep from getting into trouble. We imagine a modification of the experiment of Fig. 1-3, in which the wall with the holes consists of a plate mounted on rollers so that it can move freely up and down (in the x -direction), as shown in Fig. 1-6. By watching the motion of the plate carefully we can try to tell which hole an electron goes through. Imagine what happens when the detector is placed at $x = 0$. We would expect that an electron which passes through hole 1 must be deflected downward by the plate to reach the detector. Since the vertical component of the electron momentum is changed, the plate must recoil with an equal momentum in the opposite direction. The plate will get an upward kick. If the electron goes through the lower hole, the plate should feel a downward kick. It is clear that for every position of the detector, the momentum received by the plate will have a different value for a traversal via hole 1 than for a traversal via hole 2. So! Without disturbing the electrons *at all*, but just by watching the *plate*, we can tell which path the electron used.

Now in order to do this it is necessary to know what the momentum of the screen is, before the electron goes through. So when we measure the momentum after the electron goes by, we can figure out how much the plate’s momentum has changed. But remember, according to the uncertainty principle we cannot at the same time know the position of the plate with an arbitrary accuracy. But if we do not know exactly *where* the plate is, we cannot say precisely where the two holes are. They will be in a different place for every electron that goes through. This means that the center of our interference pattern will have a different location for each electron. The wiggles of the interference pattern will be smeared out. We shall show quantitatively in the next chapter that if we determine the momentum of the plate sufficiently accurately to determine from the recoil measurement which hole was used, then the uncertainty in the x -position of the plate will, according to the uncertainty principle, be enough to shift the pattern observed at the detector up and down in the x -direction about the distance from a maximum to its nearest minimum. Such a random shift is just enough to smear out the pattern so that no interference is observed.

The uncertainty principle “protects” quantum mechanics. Heisenberg recognized that if it were possible to measure the momentum and the position simultaneously with a greater accuracy, the quantum mechanics would collapse. So he proposed that it must be impossible. Then people sat down and tried to figure out ways of doing it, and nobody could figure out a way to measure the position and the momentum of anything—a screen, an electron, a billiard ball, anything—with any greater accuracy. Quantum mechanics maintains its perilous but still correct existence.

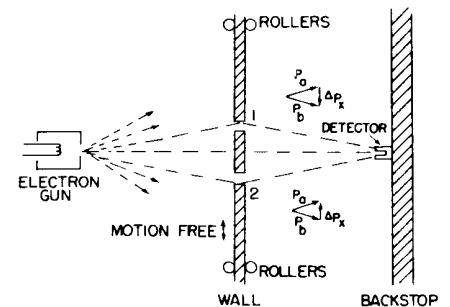


Fig. 1-6. An experiment in which the recoil of the wall is measured.

The Relation of Wave and Particle Viewpoints

2-1 Probability wave amplitudes

In this chapter we shall discuss the relationship of the wave and particle viewpoints. We already know, from the last chapter, that neither the wave viewpoint nor the particle viewpoint is correct. We would always like to present things accurately, or at least precisely enough that they will not have to be changed when we learn more—it may be extended, but it will not be changed! But when we try to talk about the wave picture or the particle picture, both are approximate, and both will change. Therefore what we learn in this chapter will not be accurate in a certain sense; we will deal with some half-intuitive arguments which will be made more precise later. But certain things will be changed a little bit when we interpret them correctly in quantum mechanics. We are doing this so that you can have some qualitative feeling for some quantum phenomena before we get into the mathematical details of quantum mechanics. Furthermore, all our experiences are with waves and with particles, and so it is rather handy to use the wave and particle ideas to get some understanding of what happens in given circumstances before we know the complete mathematics of the quantum-mechanical amplitudes. We shall try to indicate the weakest places as we go along, but most of it is very nearly correct—it is just a matter of interpretation.

First of all, we know that the new way of representing the world in quantum mechanics—the new framework—is to give an amplitude for every event that can occur, and if the event involves the reception of one particle, then we can give the amplitude to find that one particle at different places and at different times. The probability of finding the particle is then proportional to the absolute square of the amplitude. In general, the amplitude to find a particle in different places at different times varies with position and time.

In some special case it can be that the amplitude varies sinusoidally in space and time like $e^{i(\omega t - \mathbf{k} \cdot \mathbf{r})}$, where \mathbf{r} is the vector position from some origin. (Do not forget that these amplitudes are complex numbers, not real numbers.) Such an amplitude varies according to a definite frequency ω and wave number \mathbf{k} . Then it turns out that this corresponds to a classical limiting situation where we would have believed that we have a particle whose energy E was known and is related to the frequency by

$$E = \hbar\omega, \quad (2.1)$$

and whose momentum \mathbf{p} is also known and is related to the wave number by

$$\mathbf{p} = \hbar\mathbf{k}. \quad (2.2)$$

(The symbol \hbar represents the number h divided by 2π ; $\hbar = h/2\pi$.)

This means that the idea of a particle is limited. The idea of a particle—its location, its momentum, etc.—which we use so much, is in certain ways unsatisfactory. For instance, if an amplitude to find a particle at different places is given by $e^{i(\omega t - \mathbf{k} \cdot \mathbf{r})}$, whose absolute square is a constant, that would mean that the probability of finding a particle is the same at all points. That means we do not know *where* it is—it can be anywhere—there is a great uncertainty in its location.

On the other hand, if the position of a particle is more or less well known and we can predict it fairly accurately, then the probability of finding it in different places must be confined to a certain region, whose length we call Δx . Outside this region, the probability is zero. Now this probability is the absolute square of an amplitude, and if the absolute square is zero, the amplitude is also zero, so that

2-1 Probability wave amplitudes

2-2 Measurement of position and momentum

2-3 Crystal diffraction

2-4 The size of an atom

2-5 Energy levels

2-6 Philosophical implications

Note: This chapter is almost exactly the same as Chapter 38 of Volume I.

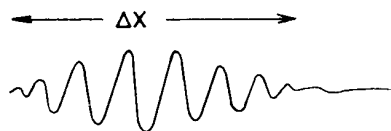


Fig. 2-1. A wave packet of length Δx .

we have a wave train whose length is Δx (Fig. 2-1), and the wavelength (the distance between nodes of the waves in the train) of that wave train is what corresponds to the particle momentum.

Here we encounter a strange thing about waves; a very simple thing which has nothing to do with quantum mechanics strictly. It is something that anybody who works with waves, even if he knows no quantum mechanics, knows: namely, *we cannot define a unique wavelength for a short wave train*. Such a wave train does not *have* a definite wavelength; there is an indefiniteness in the wave number that is related to the finite length of the train, and thus there is an indefiniteness in the momentum.

2-2 Measurement of position and momentum

Let us consider two examples of this idea—to see the reason that there is an uncertainty in the position and/or the momentum, if quantum mechanics is right. We have also seen before that if there were not such a thing—if it were possible to measure the position and the momentum of anything simultaneously—we would have a paradox; it is fortunate that we do not have such a paradox, and the fact that such an uncertainty comes naturally from the wave picture shows that everything is mutually consistent.

Here is one example which shows the relationship between the position and the momentum in a circumstance that is easy to understand. Suppose we have a single slit, and particles are coming from very far away with a certain energy—so that they are all coming essentially horizontally (Fig. 2-2). We are going to concentrate on the vertical components of momentum. All of these particles have a certain horizontal momentum p_0 , say, in a classical sense. So, in the classical sense, the vertical momentum p_y , before the particle goes through the hole, is definitely known. The particle is moving neither up nor down, because it came from a source that is far away—and so the vertical momentum is of course zero. But now let us suppose that it goes through a hole whose width is B . Then after it has come out through the hole, we know the position vertically—the y -position—with considerable accuracy—namely $\pm B/2$.[†] That is, the uncertainty in position, Δy , is of order B . Now we might also want to say, since we known the momentum is absolutely horizontal, that Δp_y is zero; but that is wrong. We *once* knew the momentum was horizontal, but we do not know it any more. Before the particles passed through the hole, we did not know their vertical positions. Now that we have found the vertical position by having the particle come through the hole, we have lost our information on the vertical momentum! Why? According to the wave theory, there is a spreading out, or diffraction, of the waves after they go through the slit, just as for light. Therefore there is a certain probability that particles coming out of the slit are not coming exactly straight. The pattern is spread out by the diffraction effect, and the angle of spread, which we can define as the angle of the first minimum, is a measure of the uncertainty in the final angle.

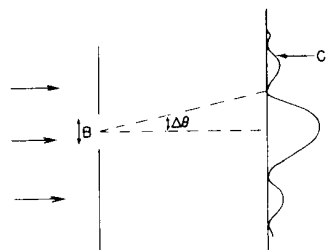


Fig. 2-2. Diffraction of particles passing through a slit.

How does the pattern become spread? To say it is spread means that there is some chance for the particle to be moving up or down, that is, to have a component of momentum up or down. We say *chance* and *particle* because we can detect this diffraction pattern with a particle counter, and when the counter receives the particle, say at C in Fig. 2-2, it receives the *entire* particle, so that, in a classical sense, the particle has a vertical momentum, in order to get from the slit up to C .

To get a rough idea of the spread of the momentum, the vertical momentum p_y has a spread which is equal to $p_0 \Delta\theta$, where p_0 is the horizontal momentum. And how big is $\Delta\theta$ in the spread-out pattern? We know that the first minimum occurs at an angle $\Delta\theta$ such that the waves from one edge of the slit have to travel one wavelength farther than the waves from the other side—we worked that out before (Chapter 30 of Vol. I). Therefore $\Delta\theta$ is λ/B , and so Δp_y in this experiment is $p_0 \lambda/B$. Note that if we make B smaller and make a more accurate measurement

[†] More precisely, the error in our knowledge of y is $\pm B/2$. But we are now only interested in the general idea, so we won't worry about factors of 2.

of the position of the particle, the diffraction pattern gets wider. So the narrower we make the slit, the wider the pattern gets, and the more is the likelihood that we would find that the particle has sidewise momentum. Thus the uncertainty in the vertical momentum is inversely proportional to the uncertainty of y . In fact, we see that the product of the two is equal to $p_0\lambda$. But λ is the wavelength and p_0 is the momentum, and in accordance with quantum mechanics, the wavelength times the momentum is Planck's constant h . So we obtain the rule that the uncertainties in the vertical momentum and in the vertical position have a product of the order h :

$$\Delta y \Delta p_y \approx h. \quad (2.3)$$

We cannot prepare a system in which we know the vertical position of a particle and can predict how it will move vertically with greater certainty than given by (2.3). That is, the uncertainty in the vertical momentum must exceed $h/\Delta y$, where Δy is the uncertainty in our knowledge of the position.

Sometimes people say quantum mechanics is all wrong. When the particle arrived from the left, its vertical momentum was zero. And now that it has gone through the slit, its position is known. Both position and momentum seem to be known with arbitrary accuracy. It is quite true that we can receive a particle, and on reception determine what its position is and what its momentum would have had to have been to have gotten there. That is true, but that is not what the uncertainty relation (2.3) refers to. Equation (2.3) refers to the *predictability* of a situation, not remarks about the *past*. It does no good to say "I knew what the momentum was before it went through the slit, and now I know the position," because now the momentum knowledge is lost. The fact that it went through the slit no longer permits us to predict the vertical momentum. We are talking about a predictive theory, not just measurements after the fact. So we must talk about what we can predict.

Now let us take the thing the other way around. Let us take another example of the same phenomenon, a little more quantitatively. In the previous example we measured the momentum by a classical method. Namely, we considered the direction and the velocity and the angles, etc., so we got the momentum by classical analysis. But since momentum is related to wave number, there exists in nature still another way to measure the momentum of a particle—photon or otherwise—which has no classical analog, because it uses Eq. (2.2). We measure the *wavelengths of the waves*. Let us try to measure momentum in this way.

Suppose we have a grating with a large number of lines (Fig. 2-3), and send a beam of particles at the grating. We have often discussed this problem: if the particles have a definite momentum, then we get a very sharp pattern in a certain direction, because of the interference. And we have also talked about how accurately we can determine that momentum, that is to say, what the resolving power of such a grating is. Rather than derive it again, we refer to Chapter 30 of Volume I, where we found that the relative uncertainty in the wavelength that can be measured with a given grating is $1/Nm$, where N is the number of lines on the grating and m is the order of the diffraction pattern. That is,

$$\Delta\lambda/\lambda = 1/Nm. \quad (2.4)$$

Now formula (2.4) can be rewritten as

$$\Delta\lambda/\lambda^2 = 1/Nm\lambda = 1/L, \quad (2.5)$$

where L is the distance shown in Fig. 2-3. This distance is the difference between the total distance that the particle or wave or whatever it is has to travel if it is reflected from the bottom of the grating, and the distance that it has to travel if it is reflected from the top of the grating. That is, the waves which form the diffraction pattern are waves which come from different parts of the grating. The first ones that arrive come from the bottom end of the grating, from the beginning of the wave train, and the rest of them come from later parts of the wave train, coming from different parts of the grating, until the last one finally arrives, and that involves a point in the wave train a distance L behind the first point. So in order that we

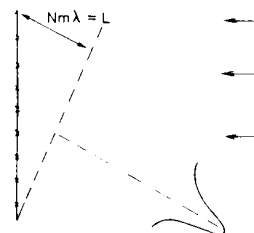


Fig. 2-3. Determination of momentum by using a diffraction grating.

shall have a sharp line in our spectrum corresponding to a definite momentum, with an uncertainty given by (2.4), we have to have a wave train of at least length L . If the wave train is too short, we are not using the entire grating. The waves which form the spectrum are being reflected from only a very short sector of the grating if the wave train is too short, and the grating will not work right—we will find a big angular spread. In order to get a narrower one, we need to use the whole grating, so that at least at some moment the whole wave train is scattering simultaneously from all parts of the grating. Thus the wave train must be of length L in order to have an uncertainty in the wavelength less than that given by (2.5). Incidentally,

$$\Delta\lambda/\lambda^2 = \Delta(1/\lambda) = \Delta k/2\pi. \quad (2.6)$$

Therefore

$$\Delta k = 2\pi/L, \quad (2.7)$$

where L is the length of the wave train.

This means that if we have a wave train whose length is less than L , the uncertainty in the wave number must exceed $2\pi/L$. Or the uncertainty in a wave number times the length of the wave train—we will call that for a moment Δx —exceeds 2π . We call it Δx because that is the uncertainty in the location of the particle. If the wave train exists only in a finite length, then that is where we could find the particle, within an uncertainty Δx . Now this property of waves, that the length of the wave train times the uncertainty of the wave number associated with it is at least 2π , is a property that is known to everyone who studies them. It has nothing to do with quantum mechanics. It is simply that if we have a finite train, we cannot count the waves in it very precisely.

Let us try another way to see the reason for that. Suppose that we have a finite train of length L ; then because of the way it has to decrease at the ends, as in Fig. 2-1, the number of waves in the length L is uncertain by something like ± 1 . But the number of waves in L is $kL/2\pi$. Thus k is uncertain, and we again get the result (2.7), a property merely of waves. The same thing works whether the waves are in space and k is the number of radians per centimeter and L is the length of the train, or the waves are in time and ω is the number of oscillations per second and T is the “length” in time that the wave train comes in. That is, if we have a wave train lasting only for a certain finite time T , then the uncertainty in the frequency is given by

$$\Delta\omega = 2\pi/T. \quad (2.8)$$

We have tried to emphasize that these are properties of waves alone, and they are well known, for example, in the theory of sound.

The point is that in quantum mechanics we interpret the wave number as being a measure of the momentum of a particle, with the rule that $p = \hbar k$, so that relation (2.7) tells us that $\Delta p \approx \hbar/\Delta x$. This, then, is a limitation of the classical idea of momentum. (Naturally, it has to be limited in some ways if we are going to represent particles by waves!) It is nice that we have found a rule that gives us some idea of when there is a failure of classical ideas.

2-3 Crystal diffraction

Next let us consider the reflection of particle waves from a crystal. A crystal is a thick thing which has a whole lot of similar atoms—we will include some complications later—in a nice array. The question is how to set the array so that we get a strong reflected maximum in a given direction for a given beam of, say, light (x-rays), electrons, neutrons, or anything else. In order to obtain a strong reflection, the scattering from all of the atoms must be in phase. There cannot be equal numbers in phase and out of phase, or the waves will cancel out. The way to arrange things is to find the regions of constant phase, as we have already explained; they are planes which make equal angles with the initial and final directions (Fig. 2-4).

If we consider two parallel planes, as in Fig. 2-4, the waves scattered from the two planes will be in phase, provided the difference in distance traveled by a wave

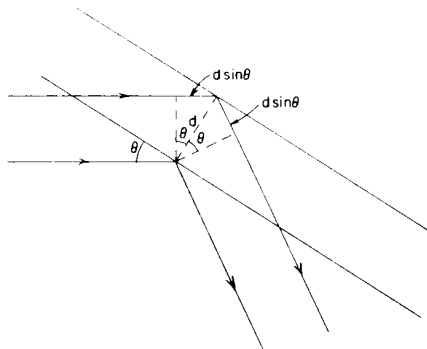


Fig. 2-4. Scattering of waves by crystal planes.

front is an integral number of wavelengths. This difference can be seen to be $2d \sin \theta$, where d is the perpendicular distance between the planes. Thus the condition for coherent reflection is

$$2d \sin \theta = n\lambda \quad (n = 1, 2, \dots). \quad (2.9)$$

If, for example, the crystal is such that the atoms happen to lie on planes obeying condition (2.9) with $n = 1$, then there will be a strong reflection. If, on the other hand, there are other atoms of the same nature (equal in density) halfway between, then the intermediate planes will also scatter equally strongly and will interfere with the others and produce no effect. So d in (2.9) must refer to *adjacent* planes; we cannot take a plane five layers farther back and use this formula!

As a matter of interest, actual crystals are not usually as simple as a single kind of atom repeated in a certain way. Instead, if we make a two-dimensional analog, they are much like wallpaper, in which there is some kind of figure which repeats all over the wallpaper. By "figure" we mean, in the case of atoms, some arrangement—calcium and a carbon and three oxygens, etc., for calcium carbonate, and so on—which may involve a relatively large number of atoms. But whatever it is, the figure is repeated in a pattern. This basic figure is called a *unit cell*.

The basic pattern of repetition defines what we call the *lattice type*; the lattice type can be immediately determined by looking at the reflections and seeing what their symmetry is. In other words, where we find any reflections *at all* determines the lattice type, but in order to determine what is in each of the elements of the lattice one must take into account the *intensity* of the scattering at the various directions. *Which* directions scatter depends on the type of lattice, but *how strongly* each scatters is determined by what is inside each unit cell, and in that way the structure of crystals is worked out.

Two photographs of x-ray diffraction patterns are shown in Figs. 2-5 and 2-6; they illustrate scattering from rock salt and myoglobin, respectively.

Incidentally, an interesting thing happens if the spacings of the nearest planes are less than $\lambda/2$. In this case (2.9) has no solution for n . Thus if λ is bigger than twice the distance between adjacent planes, then there is no side diffraction pattern, and the light—or whatever it is—will go right through the material without bouncing off or getting lost. So in the case of light, where λ is much bigger than the spacing, of course it does go through and there is no pattern of reflection from the planes of the crystal.

This fact also has an interesting consequence in the case of piles which make neutrons (these are obviously particles, for anybody's money!). If we take these neutrons and let them into a long block of graphite, the neutrons diffuse and work their way along (Fig. 2-7). They diffuse because they are bounced by the atoms, but strictly, in the wave theory, they are bounced by the atoms because of diffraction from the crystal planes. It turns out that if we take a very long piece of graphite, the neutrons that come out the far end are all of long wavelength! In fact, if one plots the intensity as a function of wavelength, we get nothing except for wavelengths longer than a certain minimum (Fig. 2-8). In other words, we can get very slow neutrons that way. Only the slowest neutrons come through; they are not diffracted or scattered by the crystal planes of the graphite, but keep going right through like light through glass, and are not scattered out the sides. There are many other demonstrations of the reality of neutron waves and waves of other particles.

2-4 The size of an atom

We now consider another application of the uncertainty relation, Eq. (2.3). It must not be taken too seriously; the idea is right but the analysis is not very accurate. The idea has to do with the determination of the size of atoms, and the fact that, classically, the electrons would radiate light and spiral in until they settle down right on top of the nucleus. But that cannot be right quantum-mechanically because then we would know where each electron was and how fast it was moving.

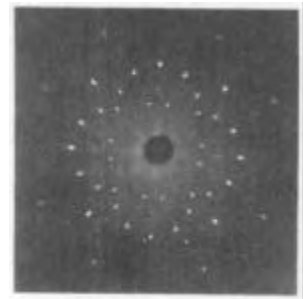


Fig. 2-5. The pattern produced by the diffraction of a beam of x-rays in a crystal of sodium chloride.

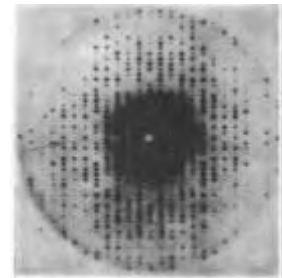


Fig. 2-6. The x-ray diffraction pattern of myoglobin.

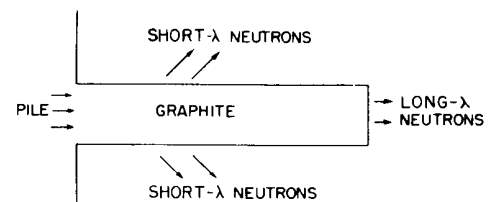


Fig. 2-7. Diffusion of pile neutrons through graphite block.

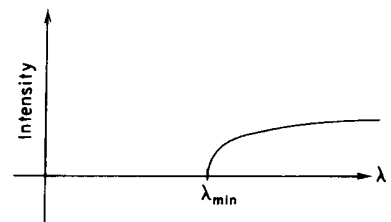


Fig. 2-8. Intensity of neutrons out of graphite rod as function of wavelength.

spectral frequencies was noted before quantum mechanics was discovered, and it is called the *Ritz combination principle*. This is again a mystery from the point of view of classical mechanics. Let us not belabor the point that classical mechanics is a failure in the atomic domain; we seem to have demonstrated that pretty well.

We have already talked about quantum mechanics as being represented by amplitudes which behave like waves, with certain frequencies and wave numbers. Let us observe how it comes about from the point of view of amplitudes that the atom has definite energy states. This is something we cannot understand from what has been said so far, but we are all familiar with the fact that confined waves have definite frequencies. For instance, if sound is confined to an organ pipe, or anything like that, then there is more than one way that the sound can vibrate, but for each such way there is a definite frequency. Thus an object in which the waves are confined has certain resonance frequencies. It is therefore a property of waves in a confined space—a subject which we will discuss in detail with formulas later on—that they exist only at definite frequencies. And since the general relation exists between frequencies of the amplitude and energy, we are not surprised to find definite energies associated with electrons bound in atoms.

2-6 Philosophical implications

Let us consider briefly some philosophical implications of quantum mechanics. As always, there are two aspects of the problem: one is the philosophical implication for physics, and the other is the extrapolation of philosophical matters to other fields. When philosophical ideas associated with science are dragged into another field, they are usually completely distorted. Therefore we shall confine our remarks as much as possible to physics itself.

First of all, the most interesting aspect is the idea of the uncertainty principle; making an observation affects the phenomenon. It has always been known that making observations affects a phenomenon, but the point is that the effect cannot be disregarded or minimized or decreased arbitrarily by rearranging the apparatus. When we look for a certain phenomenon we cannot help but disturb it in a certain minimum way, and *the disturbance is necessary for the consistency of the viewpoint*. The observer was sometimes important in prequantum physics, but only in a trivial sense. The problem has been raised: if a tree falls in a forest and there is nobody there to hear it, does it make a noise? A *real* tree falling in a *real* forest makes a sound, of course, even if nobody is there. Even if no one is present to hear it, there are other traces left. The sound will shake some leaves, and if we were careful enough we might find somewhere that some thorn had rubbed against a leaf and made a tiny scratch that could not be explained unless we assumed the leaf were vibrating. So in a certain sense we would have to admit that there is a sound made. We might ask: was there a *sensation* of sound? No, sensations have to do, presumably, with consciousness. And whether ants are conscious and whether there were ants in the forest, or whether the tree was conscious, we do not know. Let us leave the problem in that form.

Another thing that people have emphasized since quantum mechanics was developed is the idea that we should not speak about those things which we cannot measure. (Actually relativity theory also said this.) Unless a thing can be defined by measurement, it has no place in a theory. And since an accurate value of the momentum of a localized particle cannot be defined by measurement it therefore has no place in the theory. The idea that this is what was the matter with classical theory *is a false position*. It is a careless analysis of the situation. Just because we cannot *measure* position and momentum precisely does not *a priori* mean that we *cannot* talk about them. It only means that we *need* not talk about them. The situation in the sciences is this: A concept or an idea which cannot be measured or cannot be referred directly to experiment may or may not be useful. It need not exist in a theory. In other words, suppose we compare the classical theory of the world with the quantum theory of the world, and suppose that it is true experimentally that we can measure position and momentum only imprecisely. The question is whether the *ideas* of the exact position of a particle and the exact

momentum of a particle are valid or not. The classical theory admits the ideas; the quantum theory does not. This does not in itself mean that classical physics is wrong. When the new quantum mechanics was discovered, the classical people—which included everybody except Heisenberg, Schrödinger, and Born—said: “Look, your theory is not any good because you cannot answer certain questions like: what is the exact position of a particle?, which hole does it go through?, and some others.” Heisenberg’s answer was: “I do not need to answer such questions because you cannot ask such a question experimentally.” It is that we do not *have* to. Consider two theories (a) and (b); (a) contains an idea that cannot be checked directly but which is used in the analysis, and the other, (b), does not contain the idea. If they disagree in their predictions, one could not claim that (b) is false because it cannot explain this idea that is in (a), because that idea is one of the things that cannot be checked directly. It is always good to know which ideas cannot be checked directly, but it is not necessary to remove them all. It is not true that we can pursue science completely by using only those concepts which are directly subject to experiment.

In quantum mechanics itself there is a probability amplitude, there is a potential, and there are many constructs that we cannot measure directly. The basis of a science is its ability to *predict*. To predict means to tell what will happen in an experiment that has never been done. How can we do that? By assuming that we know what is there, independent of the experiment. We must extrapolate the experiments to a region where they have not been done. We must take our concepts and extend them to places where they have not yet been checked. If we do not do that, we have no prediction. So it was perfectly sensible for the classical physicists to go happily along and suppose that the position—which obviously means something for a baseball—meant something also for an electron. It was not stupidity. It was a sensible procedure. Today we say that the law of relativity is supposed to be true at all energies, but someday somebody may come along and say how stupid we were. We do not know where we are “stupid” until we “stick our neck out,” and so the whole idea is to put our neck out. And the only way to find out that we are wrong is to find out *what* our predictions are. It is absolutely necessary to make constructs.

We have already made a few remarks about the indeterminacy of quantum mechanics. That is, that we are unable now to predict what will happen in physics in a given physical circumstance which is arranged as carefully as possible. If we have an atom that is in an excited state and so is going to emit a photon, we cannot say *when* it will emit the photon. It has a certain amplitude to emit the photon at any time, and we can predict only a probability for emission; we cannot predict the future exactly. This has given rise to all kinds of nonsense and questions on the meaning of freedom of will, and of the idea that the world is uncertain.

Of course we must emphasize that classical physics is also indeterminate, in a sense. It is usually thought that this indeterminacy, that we cannot predict the future, is an important quantum-mechanical thing, and this is said to explain the behavior of the mind, feelings of free will, etc. But if the world *were* classical—if the laws of mechanics were classical—it is not quite obvious that the mind would not feel more or less the same. It is true classically that if we knew the position and the velocity of every particle in the world, or in a box of gas, we could predict exactly what would happen. And therefore the classical world is deterministic. Suppose, however, that we have a finite accuracy and do not know *exactly* where just one atom is, say to one part in a billion. Then as it goes along it hits another atom, and because we did not know the position better than to one part in a billion, we find an even larger error in the position after the collision. And that is amplified, of course, in the next collision, so that if we start with only a tiny error it rapidly magnifies to a very great uncertainty. To give an example: if water falls over a dam, it splashes. If we stand nearby, every now and then a drop will land on our nose. This appears to be completely random, yet such a behavior would be predicted by purely classical laws. The exact position of all the drops depends upon the precise wiggles of the water before it goes over the dam. How? The tiniest irregularities are magnified in falling, so that we get complete randomness. Ob-

viously, we cannot really predict the position of the drops unless we know the motion of the water *absolutely exactly*.

Speaking more precisely, given an arbitrary accuracy, no matter how precise, one can find a time long enough that we cannot make predictions valid for that long a time. Now the point is that this length of time is not very large. It is not that the time is millions of years if the accuracy is one part in a billion. The time goes, in fact, only logarithmically with the error, and it turns out that in only a very, very tiny time we lose all our information. If the accuracy is taken to be one part in billions and billions and billions—no matter how many billions we wish, provided we do stop somewhere—then we can find a time less than the time it took to state the accuracy—after which we can no longer predict what is going to happen! It is therefore not fair to say that from the apparent freedom and indeterminacy of the human mind, we should have realized that classical “deterministic” physics could not ever hope to understand it, and to welcome quantum mechanics as a release from a “completely mechanistic” universe. For already in classical mechanics there was indeterminability from a practical point of view.

Probability Amplitudes

3-1 The laws for combining amplitudes

When Schrödinger first discovered the correct laws of quantum mechanics, he wrote an equation which described the amplitude to find a particle in various places. This equation was very similar to the equations that were already known to classical physicists—equations that they had used in describing the motion of air in a sound wave, the transmission of light, and so on. So most of the time at the beginning of quantum mechanics was spent in solving this equation. But at the same time an understanding was being developed, particularly by Born and Dirac, of the basically new physical ideas behind quantum mechanics. As quantum mechanics developed further, it turned out that there were a large number of things which were not directly encompassed in the Schrödinger equation—such as the spin of the electron, and various relativistic phenomena. Traditionally, all courses in quantum mechanics have begun in the same way, retracing the path followed in the historical development of the subject. One first learns a great deal about classical mechanics so that he will be able to understand how to solve the Schrödinger equation. Then he spends a long time working out various solutions. Only after a detailed study of this equation does he get to the “advanced” subject of the electron’s spin.

We had also originally considered that the right way to conclude these lectures on physics was to show how to solve the equations of classical physics in complicated situations—such as the description of sound waves in enclosed regions, modes of electromagnetic radiation in cylindrical cavities, and so on. That was the original plan for this course. However, we have decided to abandon that plan and to give instead an introduction to the quantum mechanics. We have come to the conclusion that what are usually called the advanced parts of quantum mechanics are, in fact, quite simple. The mathematics that is involved is particularly simple, involving simple algebraic operations and no differential equations or at most only very simple ones. The only problem is that we must jump the gap of no longer being able to describe the behavior *in detail* of particles in space. So this is what we are going to try to do: to tell you about what conventionally would be called the “advanced” parts of quantum mechanics. But they are, we assure you, by all odds the simplest parts—in a deep sense of the word—as well as the most basic parts. This is frankly a pedagogical experiment; it has never been done before, as far as we know.

In this subject we have, of course, the difficulty that the quantum mechanical behavior of things is quite strange. Nobody has an everyday experience to lean on to get a rough, intuitive idea of what will happen. So there are two ways of presenting the subject: We could either describe what can happen in a rather rough physical way, telling you more or less what happens without giving the precise laws of everything; or we could, on the other hand, give the precise laws in their abstract form. But, then because of the abstractions, you wouldn’t know what they were all about, physically. The latter method is unsatisfactory because it is completely abstract, and the first way leaves an uncomfortable feeling because one doesn’t know exactly what is true and what is false. We are not sure how to overcome this difficulty. You will notice, in fact, that Chapters 1 and 2 showed this problem. The first chapter was relatively precise; but the second chapter was a rough description of the characteristics of different phenomena. Here, we will try to find a happy medium between the two extremes.

3-1 The laws for combining amplitudes

3-2 The two-slit interference pattern

3-3 Scattering from a crystal

3-4 Identical particles

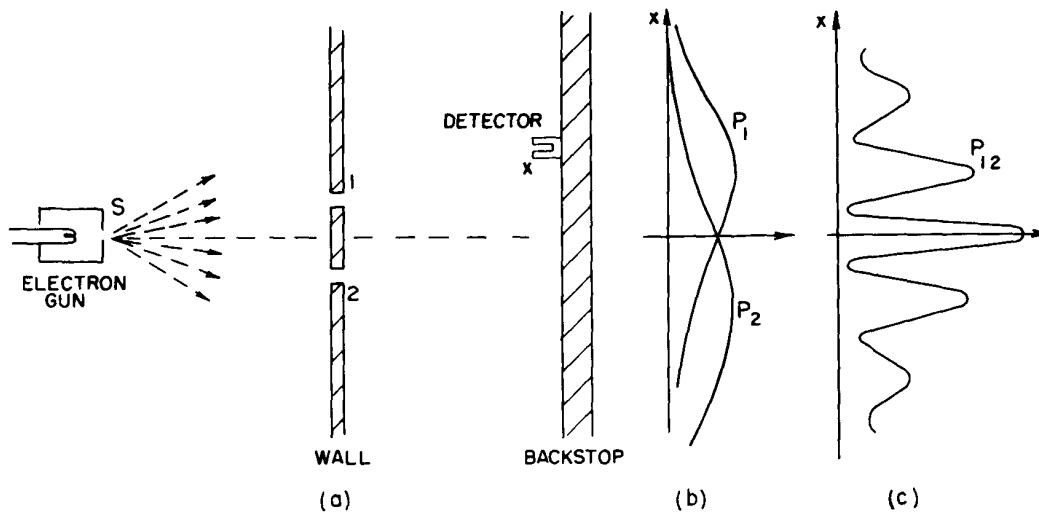


Fig 3-1. Interference experiment with electrons.

We will begin in this chapter by dealing with some general quantum mechanical ideas. Some of the statements will be quite precise, others only partially precise. It will be hard to tell you as we go along which is which, but by the time you have finished the rest of the book, you will understand in looking back which parts hold up and which parts were only explained roughly. The chapters which follow this one will not be so imprecise. In fact, one of the reasons we have tried carefully to be precise in the succeeding chapters is so that we can show you one of the most beautiful things about quantum mechanics—how much can be deduced from so little.

We begin by discussing again the superposition of *probability amplitudes*. As an example we will refer to the experiment described in Chapter 1, and shown again here in Fig. 3-1. There is a source s of particles, say electrons; then there is a wall with two slits in it; after the wall, there is a detector located at some position x . We ask for the probability that a particle will be found at x . Our *first general principle* in quantum mechanics is that the *probability* that a particle will arrive at x , when let out at the source s , can be represented quantitatively by the absolute square of a complex number called a *probability amplitude*—in this case, the “amplitude that a particle from s will arrive at x .” We will use such amplitudes so frequently that we will use a shorthand notation—invented by Dirac and generally used in quantum mechanics—to represent this idea. We write the probability amplitude this way:

$$\langle \text{Particle arrives at } x \mid \text{particle leaves } s \rangle. \quad (3.1)$$

In other words, the two brackets $\langle \rangle$ are a sign equivalent to “the amplitude that”; the expression at the *right* of the vertical line always gives the *starting* condition, and the one at the left, the *final* condition. Sometimes it will also be convenient to abbreviate still more and describe the initial and final conditions by single letters. For example, we may on occasion write the amplitude (3.1) as

$$\langle x \mid s \rangle. \quad (3.2)$$

We want to emphasize that such an amplitude is, of course, just a single number—a *complex* number.

We have already seen in the discussion of Chapter 1 that when there are two ways for the particle to reach the detector, the resulting probability is not the sum of the two probabilities, but must be written as the absolute square of the sum of two amplitudes. We had that the probability that an electron arrives at the detector when both paths are open is

$$P_{12} = |\phi_1 + \phi_2|^2. \quad (3.3)$$

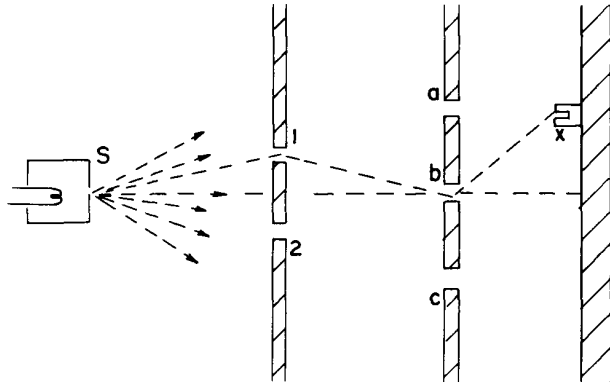


Fig. 3-2. A more complicated interference experiment.

We wish now to put this result in terms of our new notation. First, however, we want to state our *second general principle* of quantum mechanics: When a particle can reach a given state by two possible routes, the total amplitude for the process is the *sum of the amplitudes* for the two routes considered separately. In our new notation we write that

$$\langle x | s \rangle_{\text{both holes open}} = \langle x | s \rangle_{\text{through 1}} + \langle x | s \rangle_{\text{through 2}}. \quad (3.4)$$

Incidentally, we are going to suppose that the holes 1 and 2 are small enough that when we say an electron goes through the hole, we don't have to discuss which part of the hole. We could, of course, split each hole into pieces with a certain amplitude that the electron goes to the top of the hole and the bottom of the hole and so on. We will suppose that the hole is small enough so that we don't have to worry about this detail. That is part of the roughness involved; the matter can be made more precise, but we don't want to do so at this stage.

Now we want to write out in more detail what we can say about the amplitude for the process in which the electron reaches the detector at x by way of hole 1. We can do that by using our *third general principle*: When a particle goes by some particular route the amplitude for that route can be written as the *product* of the *amplitude* to go part way with the *amplitude* to go the rest of the way. For the setup of Fig. 3-1 the amplitude to go from s to x by way of hole 1 is equal to the amplitude to go from s to 1, multiplied by the amplitude to go from 1 to x .

$$\langle x | s \rangle_{\text{via 1}} = \langle x | 1 \rangle \langle 1 | s \rangle. \quad (3.5)$$

Again this result is not completely precise. We should also include a factor for the amplitude that the electron will get through the hole at 1; but in the present case it is a simple hole, and we will take this factor to be unity.

You will note that Eq. (3.5) appears to be written in reverse order. It is to be read from right to left: The electron goes from s to 1 and then from 1 to x . In summary, if events occur in succession—that is, if you can analyze one of the routes of the particle by saying it does this, then it does this, then it does that—the resultant amplitude for that route is calculated by multiplying in succession the amplitude for each of the successive events. Using this law we can rewrite Eq. (3.4) as

$$\langle x | s \rangle_{\text{both}} = \langle x | 1 \rangle \langle 1 | s \rangle + \langle x | 2 \rangle \langle 2 | s \rangle.$$

Now we wish to show that just using these principles we can calculate a much more complicated problem like the one shown in Fig. 3-2. Here we have two walls, one with two holes, 1 and 2, and another which has three holes, a , b , and c . Behind the second wall there is a detector at x , and we want to know the amplitude for a particle to arrive there. Well, one way you can find this is by calculating the superposition, or interference, of the waves that go through; but you can also do it by saying that there are six possible routes and superposing an amplitude for each. The electron can go through hole 1, then through hole a , and then to x ; or it could go through hole 1, then through hole b , and then to x ; and so on. According to our second principle, the amplitudes for alternative routes add, so we should

be able to write the amplitude from s to x as a sum of six separate amplitudes. On the other hand, using the third principle, each of these separate amplitudes can be written as a product of three amplitudes. For example, one of them is the amplitude for s to 1, times the amplitude for 1 to a , times the amplitude for a to x . Using our shorthand notation, we can write the complete amplitude to go from s to x as

$$\langle x | s \rangle = \langle x | a \rangle \langle a | 1 \rangle \langle 1 | s \rangle + \langle x | b \rangle \langle b | 1 \rangle \langle 1 | s \rangle + \cdots + \langle x | c \rangle \langle c | 2 \rangle \langle 2 | s \rangle.$$

We can save writing by using the summation notation

$$\langle x | s \rangle = \sum_{\substack{i=1,2 \\ \alpha=a,b,c}} \langle x | \alpha \rangle \langle \alpha | i \rangle \langle i | s \rangle. \quad (3.6)$$

In order to make any calculations using these methods, it is, naturally, necessary to know the amplitude to get from one place to another. We will give a rough idea of a typical amplitude. It leaves out certain things like the polarization of light or the spin of the electron, but aside from such features it is quite accurate. We give it so that you can solve problems involving various combinations of slits. Suppose a particle with a definite energy is going in empty space from a location \mathbf{r}_1 to a location \mathbf{r}_2 . In other words, it is a free particle with no forces on it. Except for a numerical factor in front, the amplitude to go from \mathbf{r}_1 to \mathbf{r}_2 is

$$\langle \mathbf{r}_2 | \mathbf{r}_1 \rangle = \frac{e^{i\mathbf{p} \cdot \mathbf{r}_{12} / \hbar}}{r_{12}}, \quad (3.7)$$

where $\mathbf{r}_{12} = \mathbf{r}_2 - \mathbf{r}_1$, and \mathbf{p} is the momentum which is related to the energy E by the relativistic equation

$$p^2 c^2 = E^2 - (m_0 c^2)^2,$$

or the nonrelativistic equation

$$\frac{p^2}{2m} = \text{Kinetic energy}.$$

Equation (3.7) says in effect that the particle has wavelike properties, the amplitude propagating as a wave with a wave number equal to the momentum divided by \hbar .

In the most general case, the amplitude and the corresponding probability will also involve the time. For most of these initial discussions we will suppose that the source always emits the particles with a given energy so we will not need to worry about the time. But we could, in the general case, be interested in some other questions. Suppose that a particle is liberated at a certain place P at a certain time, and you would like to know the amplitude for it to arrive at some location, say \mathbf{r} , at some later time. This could be represented symbolically as the amplitude $\langle \mathbf{r}, t = t_1 | P, t = 0 \rangle$. Clearly, this will depend upon both \mathbf{r} and t . You will get different results if you put the detector in different places and measure at different times. This function of \mathbf{r} and t , in general, satisfies a differential equation which is a wave equation. For example, in a nonrelativistic case it is the Schrödinger equation. One has then a wave equation analogous to the equation for electromagnetic waves or waves of sound in a gas. However, it must be emphasized that the wave function that satisfies the equation is not like a real wave in space; one cannot picture any kind of reality to this wave as one does for a sound wave.

Although one may be tempted to think in terms of "particle waves" when dealing with one particle, it is not a good idea, for if there are, say, two particles, the amplitude to find one at \mathbf{r}_1 and the other at \mathbf{r}_2 is not a simple wave in three-dimensional space, but depends on the six space variables \mathbf{r}_1 and \mathbf{r}_2 . If we are, for example, dealing with two (or more) particles, we will need the following additional principle: Provided that the two particles do not interact, the amplitude that one particle will do one thing *and* the other one something else is the product of the two amplitudes that the two particles would do the two things separately. For example, if $\langle a | s_1 \rangle$ is the amplitude for particle 1 to go from s_1 to a , and $\langle b | s_2 \rangle$

is the amplitude for particle 2 to go from s_2 to b , the amplitude that *both* things will happen together is

$$\langle a | s_1 \rangle \langle b | s_2 \rangle.$$

There is one more point to emphasize. Suppose that we didn't know where the particles in Fig. 3-2 come from before arriving at holes 1 and 2 of the first wall. We can still make a prediction of what will happen beyond the wall (for example, the amplitude to arrive at x) provided that we are given two numbers: the amplitude to have arrived at 1 and the amplitude to have arrived at 2. In other words, because of the fact that the amplitude for successive events multiplies, as shown in Eq. (3.6), all you need to know to continue the analysis is two numbers—in this particular case $\langle 1 | s \rangle$ and $\langle 2 | s \rangle$. These two complex numbers are enough to predict all the future. That is what really makes quantum mechanics easy. It turns out that in later chapters we are going to do just such a thing when we specify a starting condition in terms of two (or a few) numbers. Of course, these numbers depend upon where the source is located and possibly other details about the apparatus, but given the two numbers, we do not need to know any more about such details.

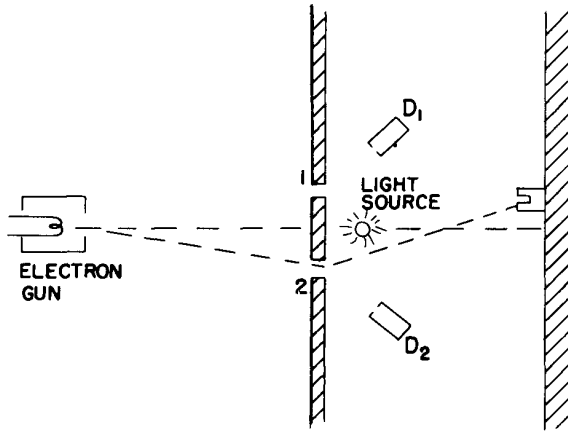


Fig. 3-3. An experiment to determine which hole the electron goes through.

3-2 The two-slit interference pattern

Now we would like to consider a matter which was discussed in some detail in Chapter 1. This time we will do it with the full glory of the amplitude idea to show you how it works out. We take the same experiment shown in Fig. 3-1, but now with the addition of a light source behind the two holes, as shown in Fig. 3-3. In Chapter 1, we discovered the following interesting result. If we looked behind slit 1 and saw a photon scattered from there, then the distribution obtained for the electrons at x in coincidence with these photons was the same as though slit 2 were closed. The total distribution for electrons that had been “seen” at either slit 1 or slit 2 was the sum of the separate distributions and was completely different from the distribution with the light turned off. This was true at least if we used light of short enough wavelength. If the wavelength was made longer so we could not be sure at which hole the scattering had occurred, the distribution became more like the one with the light turned off.

Let's examine what is happening by using our new notation and the principles of combining amplitudes. To simplify the writing, we can again let ϕ_1 stand for the amplitude that the electron will arrive at x by way of hole 1, that is,

$$\phi_1 = \langle x | 1 \rangle \langle 1 | s \rangle.$$

Similarly, we'll let ϕ_2 stand for the amplitude that the electron gets to the detector by way of hole 2:

$$\phi_2 = \langle x | 2 \rangle \langle 2 | s \rangle.$$

These are the amplitudes to go through the two holes and arrive at x if there is no light. Now if there is light, we ask ourselves the question: What is the amplitude for the process in which the electron starts at s and a photon is liberated by the

light source L , ending with the electron at x and a photon seen behind slit 1? Suppose that we observe the photon behind slit 1 by means of a detector D_1 , as shown in Fig. 3-3, and use a similar detector D_2 to count photons scattered behind hole 2. There will be an amplitude for a photon to arrive at D_1 and an electron at x , and also an amplitude for a photon to arrive at D_2 and an electron at x . Let's try to calculate them.

Although we don't have the correct mathematical formula for all the factors that go into this calculation, you will see the spirit of it in the following discussion. First, there is the amplitude $\langle 1 | s \rangle$ that an electron goes from the source to hole 1. Then we can suppose that there is a certain amplitude that while the electron is at hole 1 it scatters a photon into the detector D_1 . Let us represent this amplitude by a . Then there is the amplitude $\langle x | 1 \rangle$ that the electron goes from slit 1 to the electron detector at x . The amplitude that the electron goes from s to x via slit 1 and scatters a photon into D_1 is then

$$\langle x | 1 \rangle a \langle 1 | s \rangle.$$

Or, in our previous notation, it is just $a\phi_1$.

There is also some amplitude that an electron going through slit 2 will scatter a photon into counter D_1 . You say, "That's impossible; how can it scatter into counter D_1 if it is only looking at hole 1?" If the wavelength is long enough, there are diffraction effects, and it is certainly possible. If the apparatus is built well and if we use photons of short wavelength, then the amplitude that a photon will be scattered into detector 1, from an electron at 2 is very small. But to keep the discussion general we want to take into account that there is always some such amplitude, which we will call b . Then the amplitude that an electron goes via slit 2 and scatters a photon into D_1 is

$$\langle x | 2 \rangle b \langle 2 | s \rangle = b\phi_2.$$

The amplitude to find the electron at x and the photon in D_1 is the sum of two terms, one for each possible path for the electron. Each term is in turn made up of two factors: first, that the electron went through a hole, and second, that the photon is scattered by such an electron into detector 1; we have

$$\left\langle \begin{array}{l} \text{electron at } x \\ \text{photon at } D_1 \end{array} \middle| \begin{array}{l} \text{electron from } s \\ \text{photon from } L \end{array} \right\rangle = a\phi_1 + b\phi_2. \quad (3.8)$$

We can get a similar expression when the photon is found in the other detector D_2 . If we assume for simplicity that the system is symmetrical, then a is also the amplitude for a photon in D_2 when an electron passes through hole 2, and b is the amplitude for a photon in D_2 when the electron passes through hole 1. The corresponding total amplitude for a photon at D_2 and an electron at x is

$$\left\langle \begin{array}{l} \text{electron at } x \\ \text{photon at } D_2 \end{array} \middle| \begin{array}{l} \text{electron from } s \\ \text{photon from } L \end{array} \right\rangle = a\phi_2 + b\phi_1. \quad (3.9)$$

Now we are finished. We can easily calculate the probability for various situations. Suppose that we want to know with what probability we get a count in D_1 and an electron at x . That will be the absolute square of the amplitude given in Eq. (3.8), namely, just $|a\phi_1 + b\phi_2|^2$. Let's look more carefully at this expression. First of all, if b is zero—which is the way we would like to design the apparatus—then the answer is simply $|\phi_1|^2$ diminished in total amplitude by the factor $|a|^2$. This is the probability distribution that you would get if there were only one hole—as shown in the graph of Fig. 3-4(a). On the other hand, if the wavelength is very long, the scattering behind hole 2 into D_1 may be just about the same as for hole 1. Although there may be some phases involved in a and b , we can ask about a simple case in which the two phases are equal. If a is practically equal to b , then the total probability becomes $|\phi_1 + \phi_2|^2$ multiplied by $|a|^2$, since the common factor a can be taken out. This, however, is just the probability

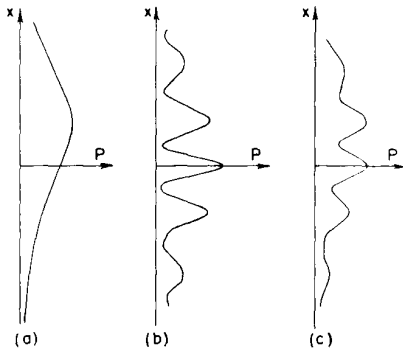


Fig. 3-4. The probability of counting an electron at x in coincidence with a photon at D in the experiment of Fig. 3-3: (a) for $b = 0$; (b) for $b = a$; (c) for $0 < b < a$.

light source L , ending with the electron at x and a photon seen behind slit 1? Suppose that we observe the photon behind slit 1 by means of a detector D_1 , as shown in Fig. 3-3, and use a similar detector D_2 to count photons scattered behind hole 2. There will be an amplitude for a photon to arrive at D_1 and an electron at x , and also an amplitude for a photon to arrive at D_2 and an electron at x . Let's try to calculate them.

Although we don't have the correct mathematical formula for all the factors that go into this calculation, you will see the spirit of it in the following discussion. First, there is the amplitude $\langle 1 | s \rangle$ that an electron goes from the source to hole 1. Then we can suppose that there is a certain amplitude that while the electron is at hole 1 it scatters a photon into the detector D_1 . Let us represent this amplitude by a . Then there is the amplitude $\langle x | 1 \rangle$ that the electron goes from slit 1 to the electron detector at x . The amplitude that the electron goes from s to x via slit 1 and scatters a photon into D_1 is then

$$\langle x | 1 \rangle a \langle 1 | s \rangle.$$

Or, in our previous notation, it is just $a\phi_1$.

There is also some amplitude that an electron going through slit 2 will scatter a photon into counter D_1 . You say, "That's impossible; how can it scatter into counter D_1 if it is only looking at hole 1?" If the wavelength is long enough, there are diffraction effects, and it is certainly possible. If the apparatus is built well and if we use photons of short wavelength, then the amplitude that a photon will be scattered into detector 1, from an electron at 2 is very small. But to keep the discussion general we want to take into account that there is always some such amplitude, which we will call b . Then the amplitude that an electron goes via slit 2 and scatters a photon into D_1 is

$$\langle x | 2 \rangle b \langle 2 | s \rangle = b\phi_2.$$

The amplitude to find the electron at x and the photon in D_1 is the sum of two terms, one for each possible path for the electron. Each term is in turn made up of two factors: first, that the electron went through a hole, and second, that the photon is scattered by such an electron into detector 1; we have

$$\left\langle \begin{array}{l} \text{electron at } x \\ \text{photon at } D_1 \end{array} \middle| \begin{array}{l} \text{electron from } s \\ \text{photon from } L \end{array} \right\rangle = a\phi_1 + b\phi_2. \quad (3.8)$$

We can get a similar expression when the photon is found in the other detector D_2 . If we assume for simplicity that the system is symmetrical, then a is also the amplitude for a photon in D_2 when an electron passes through hole 2, and b is the amplitude for a photon in D_2 when the electron passes through hole 1. The corresponding total amplitude for a photon at D_2 and an electron at x is

$$\left\langle \begin{array}{l} \text{electron at } x \\ \text{photon at } D_2 \end{array} \middle| \begin{array}{l} \text{electron from } s \\ \text{photon from } L \end{array} \right\rangle = a\phi_2 + b\phi_1. \quad (3.9)$$

Now we are finished. We can easily calculate the probability for various situations. Suppose that we want to know with what probability we get a count in D_1 and an electron at x . That will be the absolute square of the amplitude given in Eq. (3.8), namely, just $|a\phi_1 + b\phi_2|^2$. Let's look more carefully at this expression. First of all, if b is zero—which is the way we would like to design the apparatus—then the answer is simply $|\phi_1|^2$ diminished in total amplitude by the factor $|a|^2$. This is the probability distribution that you would get if there were only one hole—as shown in the graph of Fig. 3-4(a). On the other hand, if the wavelength is very long, the scattering behind hole 2 into D_1 may be just about the same as for hole 1. Although there may be some phases involved in a and b , we can ask about a simple case in which the two phases are equal. If a is practically equal to b , then the total probability becomes $|\phi_1 + \phi_2|^2$ multiplied by $|a|^2$, since the common factor a can be taken out. This, however, is just the probability

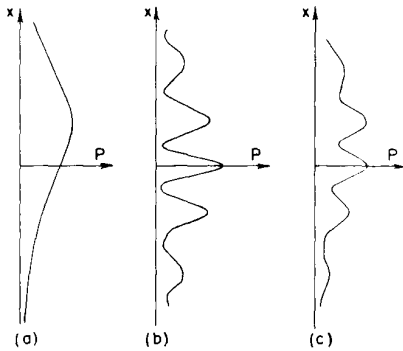


Fig. 3-4. The probability of counting an electron at x in coincidence with a photon at D in the experiment of Fig. 3-3: (a) for $b = 0$; (b) for $b = a$; (c) for $0 < b < a$.

distribution we would have gotten without the photons at all. Therefore, in the case that the wavelength is very long—and the photon detection ineffective—you return to the original distribution curve which shows interference effects, as shown in Fig. 3-4(b). In the case that the detection is partially effective, there is an interference between a lot of ϕ_1 and a little of ϕ_2 , and you will get an intermediate distribution such as is sketched in Fig. 3-4(c). Needless to say, if we look for coincidence counts of photons at D_2 and electrons at x , we will get the same kinds of results. If you remember the discussion in Chapter 1, you will see that these results give a quantitative description of what was described there.

Now we would like to emphasize an important point so that you will avoid a common error. Suppose that you only want the amplitude that the electron arrives at x , *regardless* of whether the photon was counted at D_1 or D_2 . Should you add the amplitudes given in Eqs. (3.8) and (3.9)? No! You must *never add amplitudes for different and distinct final states*. Once the photon is accepted by one of the photon counters, we can always determine which alternative occurred if we want, without any further disturbance to the system. Each alternative has a probability completely independent of the other. To repeat, do not add amplitudes for different *final* conditions, where by “final” we mean at that moment the *probability* is desired—that is, when the experiment is “finished.” You do add the amplitudes for the different *indistinguishable* alternatives inside the experiment, before the complete process is finished. At the end of the process you may say that you “don’t want to look at the photon.” That’s your business, but you still do not add the amplitudes. Nature does not know what you are looking at, and she behaves the way she is going to behave whether you bother to take down the data or not. So here we must not add the amplitudes. We first square the amplitudes for all possible different final events and then sum. The correct result for an electron at x and a photon at either D_1 or D_2 is

$$\begin{aligned} & \left| \langle e \text{ at } x \mid e \text{ from } s \rangle \right|^2 + \left| \langle e \text{ at } x \mid e \text{ from } s \rangle \right|^2 \\ & \left| \langle \text{ph at } D_1 \mid \text{ph from } L \rangle \right|^2 + \left| \langle \text{ph at } D_2 \mid \text{ph from } L \rangle \right|^2 \\ & = |a\phi_1 + b\phi_2|^2 + |a\phi_2 + b\phi_1|^2. \end{aligned} \quad (3.10)$$

3-3 Scattering from a crystal

Our next example is a phenomenon in which we have to analyze the interference of probability amplitudes somewhat carefully. We look at the process of the scattering of neutrons from a crystal. Suppose we have a crystal which has a lot of atoms with nuclei at their centers, arranged in a periodic array, and a neutron beam that comes from far away. We can label the various nuclei in the crystal by an index i , where i runs over the integers 1, 2, 3, . . . N , with N equal to the total number of atoms. The problem is to calculate the probability of getting a neutron into a counter with the arrangement shown in Fig. 3-5. For any particular atom i , the amplitude that the neutron arrives at the counter C is the amplitude that the neutron gets from the source S to nucleus i , multiplied by the amplitude a that it gets scattered there, multiplied by the amplitude that it gets from i to the counter C . Let’s write that down:

$$\langle \text{neutron at } C \mid \text{neutron from } S \rangle_{\text{via } i} = \langle C \mid i \rangle a \langle i \mid S \rangle. \quad (3.11)$$

In writing this equation we have assumed that the scattering amplitude a is the same for all atoms. We have here a large number of apparently indistinguishable routes. They are indistinguishable because a low-energy neutron is scattered from a nucleus without knocking the atom out of its place in the crystal—no “record” is left of the scattering. According to the earlier discussion, the total amplitude for a neutron at C involves a sum of Eq. (3.11) over all the atoms:

$$\langle \text{neutron at } C \mid \text{neutron from } S \rangle = \sum_{i=1}^N \langle C \mid i \rangle a \langle i \mid S \rangle. \quad (3.12)$$

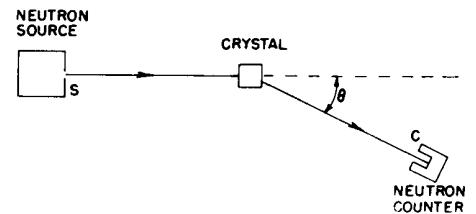


Fig. 3-5. Measuring the scattering of neutrons by a crystal.

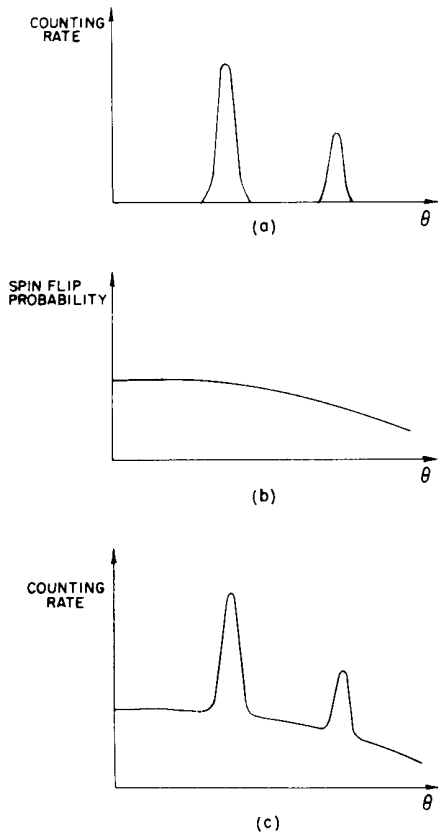


Fig. 3-6. The neutron counting rate as a function of angle: (a) for spin zero nuclei; (b) the probability of scattering with spin flip; (c) the observed counting rate with a spin one-half nucleus.

Because we are adding amplitudes of scattering from atoms with different space positions, the amplitudes will have different phases giving the characteristic interference pattern that we have already analyzed in the case of the scattering of light from a grating.

The neutron intensity as a function of angle in such an experiment is indeed often found to show tremendous variations, with very sharp interference peaks and almost nothing in between—as shown in Fig. 3-6(a). However, for certain kinds of crystals it does not work this way, and there is—along with the interference peaks discussed above—a general background of scattering in all directions. We must try to understand the apparently mysterious reasons for this. Well, we have not considered one important property of the neutron. It has a spin of one-half, and so there are two conditions in which it can be: either spin “up” (say perpendicular to the page in Fig. 3-5) or spin “down.” If the nuclei of the crystal have no spin, the neutron spin doesn’t have any effect. But when the nuclei of the crystal also have a spin, say a spin of one-half, you will observe the background of smeared-out scattering described above. The explanation is as follows.

If the neutron has one direction of spin and the atomic nucleus has the same spin, then no change of spin can occur in the scattering process. If the neutron and atomic nucleus have opposite spin, then scattering can occur by two processes, one in which the spins are unchanged and another in which the spin directions are exchanged. This rule for no net change of the sum of the spins is analogous to our classical law of conservation of angular momentum. We can begin to understand the phenomenon if we assume that all the scattering nuclei are set up with spins in one direction. A neutron with the same spin will scatter with the expected sharp interference distribution. What about one with opposite spin? If it scatters without spin flip, then nothing is changed from the above; but if the two spins flip over in the scattering, we could, in principle, find out which nucleus had done the scattering, since it would be the only one with spin turned over. Well, if we can tell which atom did the scattering, what have the other atoms got to do with it? Nothing, of course. The scattering is exactly the same as that from a single atom.

To include this effect, the mathematical formulation of Eq. (3.12) must be modified since we haven’t described the states completely in that analysis. Let’s start with all neutrons from the source having spin up and all the nuclei of the crystal having spin down. First, we would like the amplitude that at the counter the spin of the neutron is up *and* all spins of the crystal are still down. This is not different from our previous discussion. We will let a be the amplitude to scatter with no flip or spin. The amplitude for scattering from the i th atom is, of course,

$$\langle C_{\text{up}}, \text{crystal all down} | S_{\text{up}}, \text{crystal all down} \rangle = \langle C | i \rangle a \langle i | S \rangle.$$

Since all the atomic spins are still down, the various alternatives (different values of i) cannot be distinguished. There is clearly no way to tell which atom did the scattering. For this process, all the amplitudes interfere.

We have another case, however, where the spin of the detected neutron is down although it started from S with spin up. In the crystal, one of the spins must be changed to the up direction—let’s say that of the k th atom. We will assume that there is the same scattering amplitude with spin flip for every atom, namely b . (In a real crystal there is the disagreeable possibility that the reversed spin moves to some other atom, but let’s take the case of a crystal for which this probability is very low.) The scattering amplitude is then

$$\langle C_{\text{down}}, \text{nucleus } k \text{ up} | S_{\text{up}}, \text{crystal all down} \rangle = \langle C | k \rangle b \langle k | S \rangle. \quad (3.13)$$

If we ask for the probability of finding the neutron spin down and the k th nucleus spin up, it is equal to the absolute square of this amplitude, which is simply $|b|^2$ times $|\langle C | k \rangle \langle k | S \rangle|^2$. The second factor is almost independent of location in the crystal, and all phases have disappeared in taking the absolute square. The

probability of scattering from *any nucleus* in the crystal with spin flip is now

$$|b|^2 \sum_{k=1}^N |\langle C | k \rangle \langle k | S \rangle|^2,$$

which will show a smooth distribution as in Fig. 3-6(b).

You may argue, "I don't care which atom is up." Perhaps you don't, but nature knows; and the probability is, in fact, what we gave above—there isn't any interference. On the other hand, if we ask for the probability that the spin is up at the detector and all the atoms still have spin down, then we must take the absolute square of

$$\sum_{i=1}^N \langle C | i \rangle a \langle i | S \rangle.$$

Since the terms in this sum have phases, they do interfere, and we get a sharp interference pattern. If we do an experiment in which we don't observe the spin of the detected neutron, then both kinds of events can occur; and the separate probabilities add. The total probability (or counting rate) as a function of angle then looks like the graph in Fig. 3-6(c).

Let's review the physics of this experiment. If you could, *in principle*, distinguish the alternative *final* states (even though you do not bother to do so), the total, final probability is obtained by calculating the *probability* for each state (not the amplitude) and then adding them together. If you *cannot* distinguish the final states *even in principle*, then the probability amplitudes must be summed before taking the absolute square to find the actual probability. The thing you should notice particularly is that if you were to try to represent the neutron by a wave alone, you would get the same kind of distribution for the scattering of a down-spinning neutron as for an up-spinning neutron. You would have to say that the "wave" would come from all the different atoms and interfere just as for the up-spinning one with the same wavelength. But we know that is not the way it works. So as we stated earlier, we must be careful not to attribute too much reality to the waves in space. They are useful for certain problems but not for all.

3-4 Identical particles

The next experiment we will describe is one which shows one of the beautiful consequences of quantum mechanics. It again involves a physical situation in which a thing can happen in two *indistinguishable* ways, so that there is an interference of amplitudes—as is *always* true in such circumstances. We are going to discuss the scattering, at relatively low energy, of nuclei on other nuclei. We start by thinking of α -particles (which, as you know, are helium nuclei) bombarding, say, oxygen. To make it easier for us to analyze the reaction, we will look at it in the center-of-mass system, in which the oxygen nucleus and the α -particle have their velocities in opposite directions before the collision and again in exactly opposite directions after the collision. See Fig. 3-7(a). (The magnitudes of the velocities are, of course, different, since the masses are different.) We will also suppose that there is conservation of energy and that the collision energy is low enough that neither particle is broken up or left in an excited state. The reason that the two particles deflect each other is, of course, that each particle carries a positive charge and, classically speaking, there is an electrical repulsion as they go by. The scattering will happen at different angles with different probabilities, and we would like to discuss something about the angle dependence of such scatterings. (It is possible, of course, to calculate this thing classically, and it is one of the most remarkable accidents of quantum mechanics that the answer to this problem comes out the same as it does classically. This is a curious point because it happens for no other force except the inverse square law—so it is indeed an accident.)

The probability of scattering in different directions can be measured by an experiment as shown in Fig. 3-7(a). The counter at position 1 could be designed to detect only α -particles; the counter at position 2 could be designed to detect

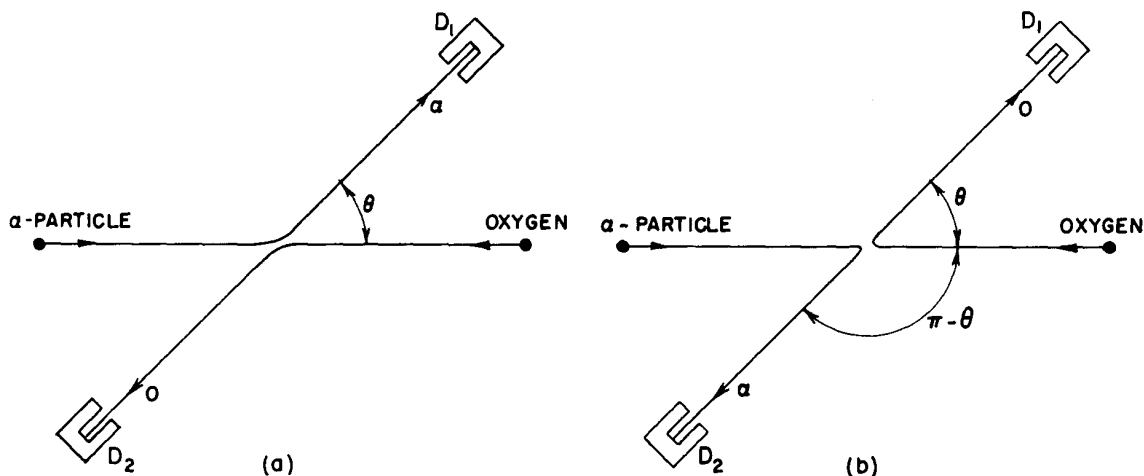


Fig. 3-7. The scattering of α -particles from oxygen nuclei, as seen in the center-of-mass system.

only oxygen—just as a check. (In the laboratory system the detectors would not be opposite; but in the CM system they are.) Our experiment consists in measuring the probability of scattering in various directions. Let's call $f(\theta)$ the amplitude to scatter into the counters when they are at the angle θ ; then $|f(\theta)|^2$ will be our experimentally determined probability.

Now we could set up another experiment in which our counters would respond to *either* the α -particle *or* the oxygen nucleus. Then we have to work out what happens when we do not bother to distinguish which particles are counted. Of course, if we are to get an oxygen in the position θ , there must be an α -particle on the opposite side at the angle $(\pi - \theta)$, as shown in Fig. 3-7(b). So if $f(\theta)$ is the amplitude for α -scattering through the angle θ , then $f(\pi - \theta)$ is the amplitude for oxygen scattering through the angle θ .[†] Thus, the probability for having *some* particle in the detector at position 1 is:

$$\text{Probability of some particle in } D_1 = |f(\theta)|^2 + |f(\pi - \theta)|^2. \quad (3.14)$$

Note that the two states are distinguishable in principle. Even though in this experiment we *do not* distinguish them, we *could*. According to the earlier discussion, then, we must add the probabilities, not the amplitudes.

The result given above is correct for a variety of target nuclei—for α -particles on oxygen, on carbon, on beryllium, on hydrogen. *But it is wrong for α -particles on α -particles.* For the one case in which both particles are exactly the same, the experimental data disagree with the prediction of (3.14). For example, the scattering probability at 90° is exactly twice what the above theory predicts and has nothing to do with the particles being “helium” nuclei. If the target is He^3 , but the projectiles are α -particles (He^4), then there is agreement. Only when the target is He^4 —so its nuclei are identical with the incoming α -particle—does the scattering vary in a peculiar way with angle.

Perhaps you can already see the explanation. There are two ways to get an α -particle into the counter: by scattering the bombarding α -particle at an angle θ , or by scattering it at an angle of $(\pi - \theta)$. How can we tell whether the bombarding particle or the target particle entered the counter? The answer is that we cannot. In the case of α -particles with α -particles there are two alternatives that cannot be distinguished. Here, we must let the probability *amplitudes* interfere by addition,

[†] In general, a scattering direction should, of course, be described by two angles, the polar angle ϕ , as well as the azimuthal angle θ . We would then say that an oxygen nucleus at (θ, ϕ) means that the α -particle is at $(\pi - \theta, \phi + \pi)$. However, for Coulomb scattering (and for many other cases), the scattering amplitude is independent of ϕ . Then the amplitude to get an oxygen at θ is the same as the amplitude to get the α -particle at $(\pi - \theta)$.

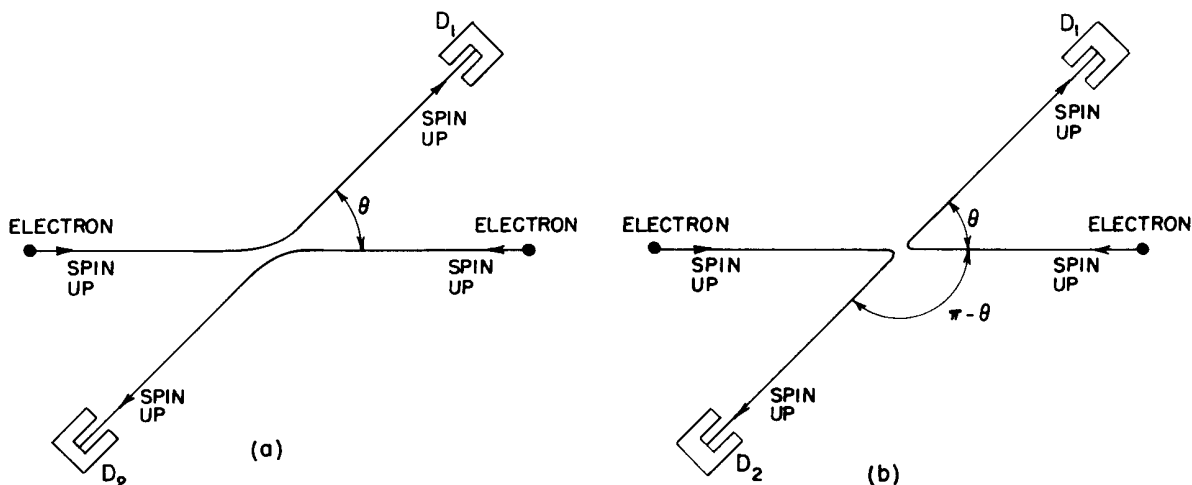


Fig. 3-8. The scattering of electrons on electrons. If the incoming electrons have parallel spins, the processes (a) and (b) are indistinguishable.

and the probability of finding an α -particle in the counter is the square of their sum:

$$\text{Probability of an } \alpha\text{-particle at } D_1 = |f(\theta) + f(\pi - \theta)|^2. \quad (3.15)$$

This is quite a different result than that in Eq. (3.14). We can take an angle of $\pi/2$ as an example, because it is easy to figure out. For $\theta = \pi/2$, we obviously have $f(\theta) = f(\pi - \theta)$, so the probability in Eq. (3.15) becomes $|f(\pi/2) + f(\pi/2)|^2 = 4|f(\pi/2)|^2$.

On the other hand, if they did not interfere, the result of Eq. (3.14) gives only $2|f(\pi/2)|^2$. So there is twice as much scattering at 90° as we might have expected. Of course, at other angles the results will also be different. And so you have the unusual result that when particles are identical, a certain new thing happens that doesn't happen when particles can be distinguished. In the mathematical description you must add the amplitudes for alternative process in which the two particles simply exchange roles and there is an interference.

An even more perplexing thing happens when we do the same kind of experiment by scattering electrons on electrons, or protons on protons. Neither of the above results is then correct! For these particles, we must invoke still a new rule, a most peculiar rule, which is the following: When you have a situation in which the identity of the electron which is arriving at a point is exchanged with another one, the new amplitude interferes with the old one with an *opposite phase*. It is interference all right, but with a minus sign. In the case of α -particles, when you exchange the α -particle entering the detector, the interfering amplitudes interfere with the positive sign. *In the case of electrons, the interfering amplitudes for exchange interfere with a negative sign.* Except for another detail to be discussed below, the proper equation for electrons in an experiment like the one shown in Fig. 3-8 is

$$\text{Probability of e at } D_1 = |f(\theta) - f(\pi - \theta)|^2. \quad (3.16)$$

The above statement must be qualified, because we have not considered the spin of the electron (α -particles have no spin). The electron spin may be considered to be either "up" or "down" with respect to the plane of the scattering. If the energy of the experiment is low enough, the magnetic forces due to the currents will be small and the spin will not be affected. We will assume that this is the case for the present analysis, so that there is no chance that the spins are changed during the collision. Whatever spin the electron has, it carries along with it. Now you see there are many possibilities. The bombarding and target particles can have both spins up, both down, or opposite spins. If both spins are up, as in Fig. 3-8 (or if both spins are down), the same will be true of the recoil particles and the *amplitude* for the process is the *difference* of the amplitudes for the two possibilities

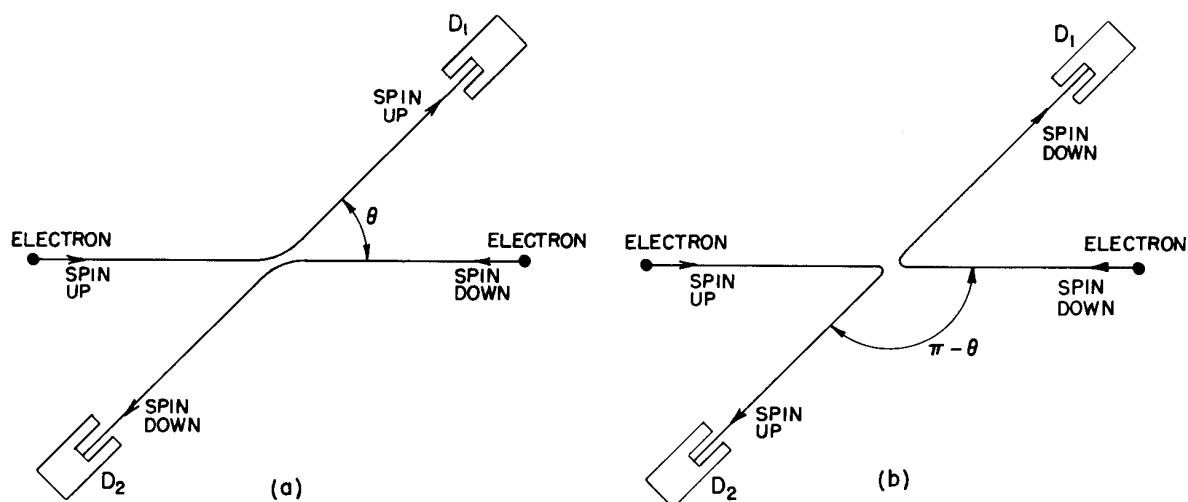


Fig. 3-9. The scattering of electrons with antiparallel spins.

shown in Fig. 3-8(a) and (b). The *probability* of detecting an electron in D_1 is then given by Eq. (3.16).

Suppose, however, the “bombarding” spin is up and the “target” spin is down. The electron entering counter 1 can have spin up or spin down, and by measuring this spin we can tell whether it came from the bombarding beam or from the target. The two possibilities are shown in Fig. 3-9(a) and (b); they are distinguishable in principle, and hence there will be no interference—merely an addition of the two probabilities. The same argument holds if both of the original spins are reversed—that is, if the left-hand spin is down and the right-hand spin is up.

Now if we take our electrons at random—as from a tungsten filament in which the electrons are completely unpolarized—then the odds are fifty-fifty that any particular electron comes out with spin up or spin down. If we don’t bother to measure the spin of the electrons at any point in the experiment, we have what we call an unpolarized experiment. The results for this experiment are best calculated by listing all of the various possibilities as we have done in Table 3-1. A separate *probability* is computed for each distinguishable alternative. The total probability is then the sum of all the separate probabilities. Note that for unpolarized beams the result for $\theta = \pi/2$ is one-half that of the classical result with independent particles. The behavior of identical particles has many interesting consequences; we will discuss them in greater detail in the next chapter.

Table 3-1

Scattering of unpolarized spin one-half particles

Fraction of cases	Spin of particle 1	Spin of particle 2	Spin at D_1	Spin at D_2	Probability
$\frac{1}{4}$	up	up	up	up	$ f(\theta) - f(\pi - \theta) ^2$
$\frac{1}{4}$	down	down	down	down	$ f(\theta) - f(\pi - \theta) ^2$
$\frac{1}{4}$	up	down	up	down	$ f(\theta) ^2$
			down	up	$ f(\pi - \theta) ^2$
$\frac{1}{4}$	down	up	up	down	$ f(\pi - \theta) ^2$
			down	up	$ f(\theta) ^2$
Total probability = $\frac{1}{2} f(\theta) - f(\pi - \theta) ^2 + \frac{1}{2} f(\theta) ^2 + \frac{1}{2} f(\pi - \theta) ^2$					

Identical Particles

4-1 Bose particles and Fermi particles

In the last chapter we began to consider the special rules for the interference that occurs in processes with two *identical* particles. By *identical particles* we mean things like electrons which can in no way be distinguished one from another. If a process involves two particles that are identical, reversing which one arrives at a counter is an alternative which cannot be distinguished and—like all cases of alternatives which cannot be distinguished—interferes with the original, un-exchanged case. The amplitude for an event is then the sum of the two interfering amplitudes; but, interestingly enough, the interference is in some cases with the *same* phase and, in others, with the *opposite* phase.

Suppose we have a collision of two particles *a* and *b* in which particle *a* scatters in the direction 1 and particle *b* scatters in the direction 2, as sketched in Fig. 4-1(a). Let's call $f(\theta)$ the amplitude for this process; then the probability P_1 of observing such an event is proportional to $|f(\theta)|^2$. Of course, it could also happen that particle *b* scattered into counter 1 and particle *a* went into counter 2, as shown in Fig. 4-1(b). Assuming that there are no special directions defined by spins or such, the probability P_2 for this process is just $|f(\pi - \theta)|^2$, because it is just equivalent to the first process with counter 1 moved over to the angle $\pi - \theta$. You might also think that the *amplitude* for the second process is just $f(\pi - \theta)$. But that is not necessarily so, because there could be an arbitrary phase factor. That is, the amplitude could be

$$e^{i\delta} f(\pi - \theta).$$

Such an amplitude still gives a probability P_2 equal to $|f(\pi - \theta)|^2$.

Now let's see what happens if *a* and *b* are identical particles. Then the two different processes shown in the two diagrams of Fig. 4-1 cannot be distinguished. There is an amplitude that *either* *a* or *b* goes into counter 1, while the other goes into counter 2. This amplitude is the sum of the amplitudes for the two processes shown in Fig. 4-1. If we call the first one $f(\theta)$, then the second one is $e^{i\delta} f(\pi - \theta)$, where now the phase factor is very important because we are going to be adding two amplitudes. Suppose we have to multiply the amplitude by a certain phase factor when we exchange the roles of the two particles. If we exchange them again we should get the same factor again. But we are then back to the first process.

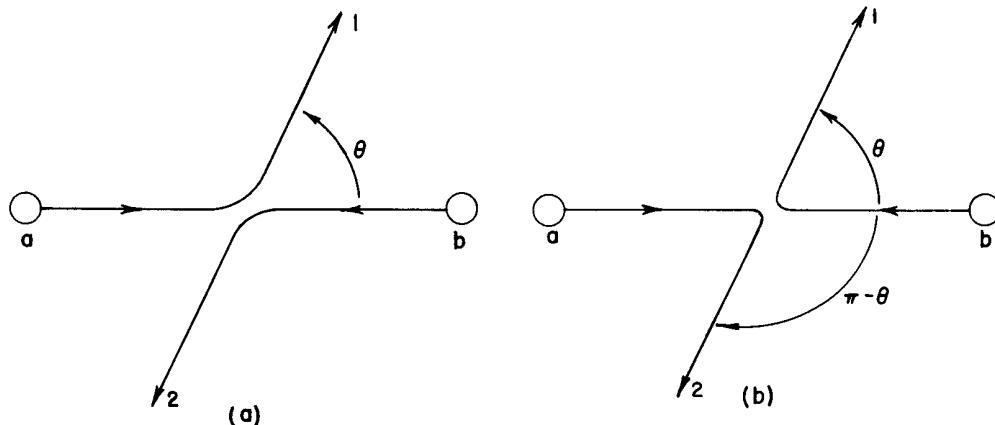


Fig. 4-1. In the scattering of two identical particles, the processes (a) and (b) are indistinguishable.

4-1 Bose particles and Fermi particles

4-2 States with two Bose particles

4-3 States with n Bose particles

4-4 Emission and absorption of photons

4-5 The blackbody spectrum

4-6 Liquid helium

4-7 The exclusion principle

Review: Blackbody radiation in:
Chapter 41, Vol. I, *The Brownian Movement*

Chapter 42, Vol. I, *Applications of Kinetic Theory*

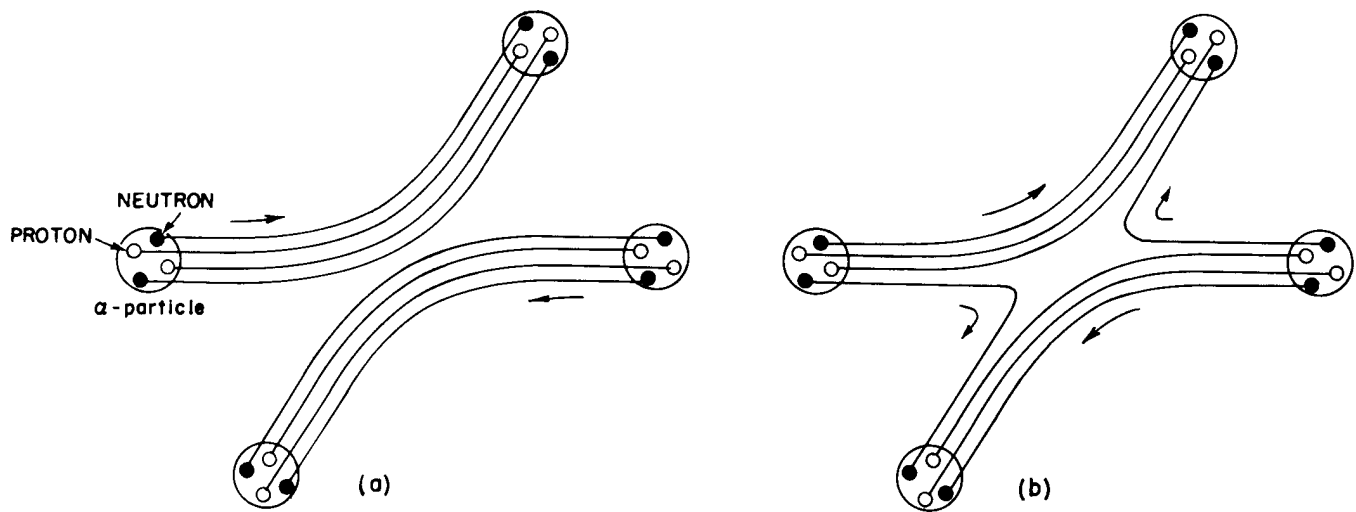


Fig. 4-2. The scattering of two α -particles. In (a) the two particles retain their identity; in (b) a neutron is exchanged during the collision.

The phase factor taken twice must bring us back where we started—its square must be equal to 1. There are only two possibilities: $e^{i\delta}$ is equal to $+1$, or is equal to -1 . Either the exchanged case contributes with the *same* sign, or it contributes with the *opposite* sign. Both cases exist in nature, each for a different class of particles. Particles which interfere with a *positive* sign are called *Bose particles* and those which interfere with a *negative* sign are called *Fermi particles*. The Bose particles are the photon, the mesons, and the graviton. The Fermi particles are the electron, the muon, the neutrinos, the nucleons, and the baryons. We have, then, that the amplitude for the scattering of identical particles is:

Bose particles:

$$(\text{Amplitude direct}) + (\text{Amplitude exchanged}). \quad (4.1)$$

Fermi particles:

$$(\text{Amplitude direct}) - (\text{Amplitude exchanged}). \quad (4.2)$$

For particles with spin—like electrons—there is an additional complication. We must specify not only the location of the particles but the direction of their spins. It is only for identical particles *with identical spin states* that the amplitudes interfere when the particles are exchanged. If you think of the scattering of unpolarized beams—which are a mixture of different spin states—there is some extra arithmetic.

Now an interesting problem arises when there are two or more particles bound tightly together. For example, an α -particle has four particles in it—two neutrons and two protons. When two α -particles scatter, there are several possibilities. It may be that during the scattering there is a certain amplitude that one of the neutrons will leap across from one α -particle to the other, while a neutron from the other α -particle leaps the other way so that the two alphas which come out of the scattering are not the original ones—there has been an exchange of a pair of neutrons. See Fig. 4-2. The amplitude for scattering with an exchange of a pair of neutrons will interfere with the amplitude for scattering with no such exchange, and the interference must be with a minus sign because there has been an exchange of one pair of Fermi particles. On the other hand, if the relative energy of the two α -particles is so low that they stay fairly far apart—say, due to the Coulomb repulsion—and there is never any appreciable probability of exchanging any of the internal particles, we can consider the α -particle as a simple object, and we do not need to worry about its internal details. In such circumstances, there are only two contributions to the scattering amplitude. Either there is no exchange, or all four of the nucleons are exchanged in the scattering. Since the protons and the

neutrons in the α -particle are all Fermi particles, an exchange of any pair reverses the sign of the scattering amplitude. So long as there are no internal changes in the α -particles, interchanging the two α -particles is the same as interchanging four pairs of Fermi particles. There is a change in sign for each pair, so the net result is that the amplitudes combine with a positive sign. The α -particle behaves like a Bose particle.

So the rule is that composite objects, in circumstances in which the composite object can be considered as a single object, behave like Fermi particles or Bose particles, depending on whether they contain an odd number or an even number of Fermi particles.

All the elementary Fermi particles we have mentioned—such as the electron, the proton, the neutron, and so on—have a spin $j = 1/2$. If several such Fermi particles are put together to form a composite object, the resulting spin may be either integral or half-integral. For example, the common isotope of helium, He^4 , which has two neutrons and two protons, has a spin of zero, whereas Li^7 which has three protons and four neutrons, has a spin of $3/2$. We will learn later the rules for compounding angular momentum, and will just mention now that every composite object which has a *half-integral spin* imitates a *Fermi particle*, whereas every composite object with an *integral spin* imitates a *Bose particle*.

This brings up an interesting question: Why is it that particles with half-integral spin are Fermi particles whose amplitudes add with the minus sign, whereas particles with integral spin are Bose particles whose amplitudes add with the positive sign? We apologize for the fact that we cannot give you an elementary explanation. An explanation has been worked out by Pauli from complicated arguments of quantum field theory and relativity. He has shown that the two must necessarily go together, but we have not been able to find a way of reproducing his arguments on an elementary level. It appears to be one of the few places in physics where there is a rule which can be stated very simply, but for which no one has found a simple and easy explanation. The explanation is deep down in relativistic quantum mechanics. This probably means that we do not have a complete understanding of the fundamental principle involved. For the moment, you will just have to take it as one of the rules of the world.

4-2 States with two Bose particles

Now we would like to discuss an interesting consequence of the addition rule for Bose particles. It has to do with their behavior when there are several particles present. We begin by considering a situation in which two Bose particles are scattered from two different scatterers. We won't worry about the details of the scattering mechanism. We are interested only in what happens to the scattered particles. Suppose we have the situation shown in Fig. 4-3. The particle a is scattered into the state 1. By a *state* we mean a given direction and energy, or some other given condition. The particle b is scattered into the state 2. We want to assume that the two states 1 and 2 are nearly the same. (What we really want to find out eventually is the amplitude that the two particles are scattered into identical directions, or states; but it is best if we think first about what happens if the states are almost the same and then work out what happens when they become identical.)

Suppose that we had only particle a ; then it would have a certain amplitude for scattering in direction 1, say $\langle 1 | a \rangle$. And particle b alone would have the amplitude $\langle 2 | b \rangle$ for landing in direction 2. If the two particles are not identical, the amplitude for the two scatterings to occur at the same time is just the product

$$\langle 1 | a \rangle \langle 2 | b \rangle.$$

The probability for such an event is then

$$|\langle 1 | a \rangle \langle 2 | b \rangle|^2,$$

which is also equal to

$$|\langle 1 | a \rangle|^2 |\langle 2 | b \rangle|^2.$$

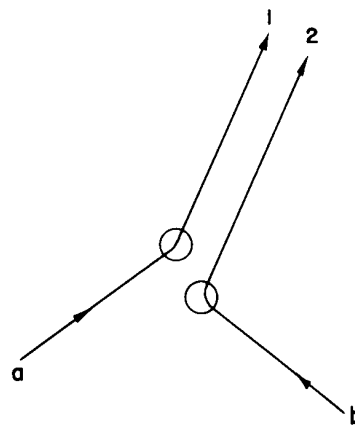


Fig. 4-3. A double scattering into nearby final states.

To save writing for the present arguments, we will sometimes set

$$\langle 1 | a \rangle = a_1, \quad \langle 2 | b \rangle = b_2.$$

Then the probability of the double scattering is

$$|a_1|^2 |b_2|^2.$$

It could also happen that particle b is scattered into direction 1, while particle a goes into direction 2. The amplitude for this process is

$$\langle 2 | a \rangle \langle 1 | b \rangle,$$

and the probability of such an event is

$$|\langle 2 | a \rangle \langle 1 | b \rangle|^2 = |a_2|^2 |b_1|^2.$$

Imagine now that we have a pair of tiny counters that pick up the two scattered particles. The probability P_2 that they will pick up two particles together is just the sum

$$P_2 = |a_1|^2 |b_2|^2 + |a_2|^2 |b_1|^2. \quad (4.3)$$

Now let's suppose that the directions 1 and 2 are very close together. We expect that a should vary smoothly with direction, so a_1 and a_2 must approach each other as 1 and 2 get close together. If they are close enough, the amplitudes a_1 and a_2 will be equal. We can set $a_1 = a_2$ and call them both just a ; similarly, we set $b_1 = b_2 = b$. Then we get that

$$P_2 = 2|a|^2 |b|^2. \quad (4.4)$$

Now suppose, however, that a and b are identical Bose particles. Then the process of a going into 1 and b going into 2 cannot be distinguished from the exchanged process in which a goes into 2 and b goes into 1. In this case the *amplitudes* for the two different processes can interfere. The *total* amplitude to obtain a particle in each of the two counters is

$$\langle 1 | a \rangle \langle 2 | b \rangle + \langle 2 | a \rangle \langle 1 | b \rangle. \quad (4.5)$$

And the probability that we get a pair is the absolute square of this amplitude,

$$P_2 = |a_1 b_2 + a_2 b_1|^2 = 4|a|^2 |b|^2. \quad (4.6)$$

We have the result that it is *twice as likely* to find two *identical* Bose particles scattered into the same state *as you would calculate assuming the particles were different*.

Although we have been considering that the two particles are observed in separate counters, this is not essential—as we can see in the following way. Let's imagine that both the directions 1 and 2 would bring the particles into a *single* small counter which is some distance away. We will let the direction 1 be defined by saying that it heads toward the element of area dS_1 of the counter. Direction 2 heads toward the surface element dS_2 of the counter. (We imagine that the counter presents a surface at right angles to the line from the scatterings.) Now we cannot give a probability that a particle will go into a precise direction or to a *particular* point in space. Such a thing is impossible—the chance for any exact direction is zero. When we want to be so specific, we shall have to define our amplitudes so that they give the probability of arriving *per unit area* of a counter. Suppose that we had only particle a ; it would have a certain amplitude for scattering in direction 1. Let's define $\langle 1 | a \rangle = a_1$ to be the amplitude that a will scatter *into a unit area* of the counter in the direction 1. In other words, the scale of a_1 is chosen—we say it is “normalized” so that the probability that it will scatter *into an element of area* dS_1 is

$$|\langle 1 | a \rangle|^2 dS_1 = |a_1|^2 dS_1. \quad (4.7)$$

If our counter has the total area ΔS , and we let dS_1 range over this area, the total probability that the particle a will be scattered into the counter is

$$\int_{\Delta S} |a_1|^2 dS_1. \quad (4.8)$$

As before, we want to assume that the counter is sufficiently small so that the amplitude a_1 doesn't vary significantly over the surface of the counter; a_1 is then a constant amplitude which we can call a . Then the probability that particle a is scattered somewhere into the counter is

$$p_a = |a|^2 \Delta S. \quad (4.9)$$

In the same way, we will have that the probability that particle b —when it is alone—scatters into some element of area, say dS_2 , is

$$|b_2|^2 dS_2.$$

(We use dS_2 instead of dS_1 because we will later want a and b to go into different directions.) Again we set b_2 equal to the constant amplitude b ; then the probability that particle b is counted in the detector is

$$p_b = |b|^2 \Delta S. \quad (4.10)$$

Now when both particles are present, the probability that a is scattered into dS_1 and b is scattered into dS_2 is

$$|a_1 b_2|^2 dS_1 dS_2 = |a|^2 |b|^2 dS_1 dS_2. \quad (4.11)$$

If we want the probability that *both* a and b get into the counter, we integrate both dS_1 and dS_2 over ΔS and find that

$$P_2 = |a|^2 |b|^2 (\Delta S)^2. \quad (4.12)$$

We notice, incidentally, that this is just equal to $p_a \cdot p_b$, just as you would suppose assuming that the particles a and b act independently of each other.

When the two particles are identical, however, there are two indistinguishable possibilities for each pair of surface elements dS_1 and dS_2 . Particle a going into dS_2 and particle b going into dS_1 is indistinguishable from a into dS_1 and b into dS_2 , so the amplitudes for these processes will interfere. (When we had two *different* particles above—although we did not *in fact* care which particle went where in the counter—we could, *in principle*, have found out; so there was no interference. For identical particles we cannot tell, even in principle.) We must write, then, that the probability that the two particles arrive at dS_1 and dS_2 is

$$|a_1 b_2 + a_2 b_1|^2 dS_1 dS_2. \quad (4.13)$$

Now, however, when we integrate over the area of the counter, we must be careful. If we let dS_1 and dS_2 range over the whole area ΔS , we would count each part of the area *twice* since (4.13) contains everything that can happen with any pair of surface elements dS_1 and dS_2 †. We can still do the integral that way, if we correct for the double counting by dividing the result by 2. We get then that P_2 for identical Bose particles is

$$P_2(\text{Bose}) = \frac{1}{2} \{4|a|^2 |b|^2 (\Delta S)^2\} = 2|a|^2 |b|^2 (\Delta S)^2. \quad (4.14)$$

Again, this is just twice what we got in Eq. (4.12) for distinguishable particles.

If we imagine for a moment that we knew that the b channel had already sent its particle into the particular direction, we can say that the *probability* that a second particle will go into the same direction is twice as great as we would have

† In (4.11) interchanging dS_1 and dS_2 gives a different event, so both surface elements should range over the whole area of the counter. In (4.13) we are treating dS_1 and dS_2 as a *pair* and including everything that can happen. If the integrals include again what happens when dS_1 and dS_2 are reversed, everything is counted twice.

expected if we had calculated it as an independent event. It is a property of Bose particles that if there is already one particle in a condition of some kind, the *probability* of getting a second one in the same condition is twice as great as it would be if the first one were not already there. This fact is often stated in the following way: If there is already one Bose particle in a given state, the amplitude for putting an identical one on top of it is $\sqrt{2}$ greater than if it weren't there. (This is not a proper way of stating the result from the physical point of view we have taken, but if it is used consistently as a rule, it will, of course, give the correct result.)

4-3 States with n Bose particles

Let's extend our result to a situation in which there are n particles present. We imagine the circumstance shown in Fig. 4-4. We have n particles a, b, c, \dots , which are scattered and end up in the directions $1, 2, 3, \dots, n$. All n directions are headed toward a small counter a long distance away. As in the last section, we choose to normalize all the amplitudes so that the probability that each particle acting alone would go into an element of surface dS of the counter is

$$|\langle \rangle|^2 dS.$$

First, let's assume that the particles are all distinguishable; then the probability that n particles will be counted together in n different surface elements is

$$|a_1 b_2 c_3 \dots|^2 dS_1 dS_2 dS_3 \dots \quad (4.15)$$

Again we take that the amplitudes don't depend on where dS is located in the counter (assumed small) and call them simply a, b, c, \dots . The probability (4.15) becomes

$$|a|^2 |b|^2 |c|^2 \dots dS_1 dS_2 dS_3 \dots \quad (4.16)$$

Integrating each dS over the surface ΔS of the counter, we have that P_n (different), the probability of counting n different particles at once, is

$$P_r \text{ (different)} = |a|^2 |b|^2 |c|^2 \dots (\Delta S)^n \quad (4.17)$$

This is just the product of the probabilities for each particle to enter the counter separately. They all act independently—the probability for one to enter does not depend on how many others are also entering.

Now suppose that all the particles are identical Bose particles. For each set of directions $1, 2, 3, \dots$ there are many indistinguishable possibilities. If there were, for instance, just three particles, we would have the following possibilities:

$a \rightarrow 1$	$a \rightarrow 1$	$a \rightarrow 2$
$b \rightarrow 2$	$b \rightarrow 3$	$b \rightarrow 1$
$c \rightarrow 3$	$c \rightarrow 2$	$c \rightarrow 3$
$a \rightarrow 2$	$a \rightarrow 3$	$a \rightarrow 3$
$b \rightarrow 3$	$b \rightarrow 1$	$b \rightarrow 2$
$c \rightarrow 1$	$c \rightarrow 2$	$c \rightarrow 1$

There are six different combinations. With n particles, there are $n!$ different, but *indistinguishable*, possibilities for which we must add amplitudes. The probability that n particles will be counted in n surface elements is then

$$|a_1 b_2 c_3 \dots + a_1 b_3 c_2 \dots + a_2 b_1 c_3 \dots + a_2 b_3 c_1 \dots + \text{etc.} + \text{etc.}|^2 dS_1 dS_2 dS_3 \dots dS_n. \quad (4.18)$$

Once more we assume that all the directions are so close that we can set $a_1 = a_2 = \dots = a = a_n$, and similarly for b, c, \dots ; the probability of (4.18) becomes

$$|n! abc \dots|^2 dS_1 dS_2 \dots dS_n. \quad (4.19)$$

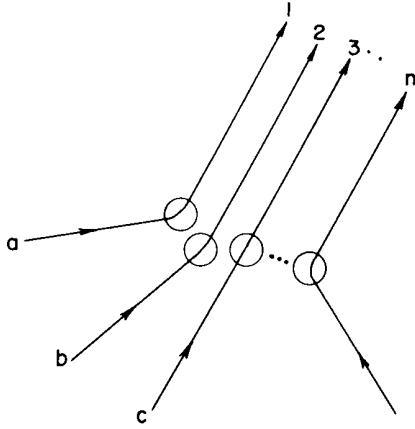


Fig. 4-4. The scattering of n particles into nearby final states.

When we integrate each dS over the area ΔS of the counter, each possible product of surface elements is counted $n!$ times; we correct for this by dividing by $n!$ and get

$$P_n(\text{Bose}) = \frac{1}{n!} |n! abc \dots|^2 (\Delta S)^n$$

or

$$P_n(\text{Bose}) = n! |abc \dots|^2 (\Delta S)^n. \quad (4.20)$$

Comparing this result with Eq. (4.17), we see that the probability of counting n Bose particles together is $n!$ greater than we would calculate assuming that the particles were all distinguishable. We can summarize our result this way:

$$P_n(\text{Bose}) = n! P_n(\text{different}). \quad (4.21)$$

Thus, the probability in the Bose case is larger by $n!$ than you would calculate assuming that the particles acted independently.

We can see better what this means if we ask the following question: What is the probability that a Bose particle will go into a particular state when there are *already* n others present? Let's call the newly added particle w . If we have $(n + 1)$ particles, including w , Eq. (4.20) becomes

$$P_{n+1}(\text{Bose}) = (n + 1)! |abc \dots w|^2 (\Delta S)^{n+1}. \quad (4.22)$$

We can write this as

$$P_{n+1}(\text{Bose}) = \{(n + 1)|w|^2 \Delta S\} n! |abc \dots|^2 \Delta S^n$$

or

$$P_{n+1}(\text{Bose}) = (n + 1)|w|^2 \Delta S P_n(\text{Bose}). \quad (4.23)$$

We can look at this result in the following way: The number $|w|^2 \Delta S$ is the probability for getting particle w into the detector if no other particles were present; $P_n(\text{Bose})$ is the chance that there are already n other Bose particles present. So Eq. (4.23) says that *when* there are n other identical Bose particles present, the probability that *one more* particle will enter the same state is *enhanced* by the factor $(n + 1)$. The probability of getting a boson, where there are already n , is $(n + 1)$ times stronger than it would be if there were none before. The *presence* of the other particles increases the probability of getting one more.

4-4 Emission and absorption of photons

Throughout our discussion we have talked about a process like the scattering of α -particles. But that is not essential; we could have been speaking of the creation of particles, as for instance the emission of light. When the light is emitted, a photon is "created." In such a case, we don't need the incoming lines in Fig. 4-4; we can consider merely that there are n atoms a, b, c, \dots emitting light, as in Fig. 4-5. So our result can also be stated: *The probability that an atom will emit a photon into a particular final state is increased by the factor $(n + 1)$ if there are already n photons in that state.*

People like to summarize this result by saying that the *amplitude* to emit a photon is increased by the factor $\sqrt{n + 1}$ when there are already n photons present. It is, of course, another way of saying the same thing if it is understood to mean that this amplitude is just to be squared to get the probability.

It is generally true in quantum mechanics that the amplitude to get from any condition ϕ to any other condition χ is the complex conjugate of the amplitude to get from χ to ϕ :

$$\langle \chi | \phi \rangle = \langle \phi | \chi \rangle^*. \quad (4.24)$$

We will learn about this law a little later, but for the moment we will just assume it is true. We can use it to find out how photons are scattered or absorbed *out* of a given state. We have that the amplitude that a photon will be added to some state, say i , when there are already n photons present is, say,

$$\langle n + 1 | i \rangle = \sqrt{n + 1} a, \quad (4.25)$$

4-7

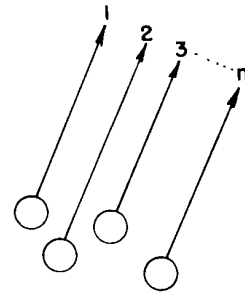


Fig. 4-5. The creation of n photons in nearby states.

where $a = \langle i | a \rangle$ is the amplitude when there are no others present. Using Eq. (4.24), the amplitude to go the other way—from $(n + 1)$ photons to n —is

$$\langle n | n + 1 \rangle = \sqrt{n + 1} a^*. \quad (4.26)$$

This isn't the way people usually say it; they don't like to think of going from $(n + 1)$ to n , but prefer always to start with n photons present. Then they say that the amplitude to absorb a photon when there are n present—in other words, to go from n to $(n + 1)$ —is

$$\langle n + 1 | n \rangle = \sqrt{n} a^*. \quad (4.27)$$

which is, of course, just the same as Eq. (4.26). Then they have trouble trying to remember when to use \sqrt{n} or $\sqrt{n + 1}$. Here's the way to remember: The factor is always the square root of the largest number of photons present, whether it is before or after the reaction. Equations (4.25) and (4.26) show that the law is really symmetric—it only appears unsymmetric if you write it as Eq. (4.27).

There are many physical consequences of these new rules; we want to describe one of them having to do with the emission of light. Suppose we imagine a situation in which photons are contained in a box—you can imagine a box with mirrors for walls. Now say that in the box we have n photons, all of the same state—the same frequency, direction, and polarization—so they can't be distinguished, and that also there is an atom in the box that can emit another photon into the same state. Then the probability that it will emit a photon is

$$(n + 1)|a|^2, \quad (4.28)$$

and the probability that it will absorb a photon is

$$n|a|^2, \quad (4.29)$$

where $|a|^2$ is the probability it would emit if no photons were present. We have already discussed these rules in a somewhat different way in Chapter 42 of Vol. I. Equation (4.29) says that the probability that an atom will *absorb* a photon and make a transition to a higher energy state is proportional to the intensity of the light shining on it. But, as Einstein first pointed out, the rate at which an atom will make a transition *downward* has two parts. There is the probability that it will make a spontaneous transition $|a|^2$, plus the probability of an induced transition $n|a|^2$, which is proportional to the intensity of the light—that is, to the number of photons present. Furthermore, as Einstein said, the coefficients of absorption and of induced emission are equal and are related to the probability of spontaneous emission. What we learn here is that if the light intensity is measured in terms of the number of photons present (instead of as the energy per unit area, and per sec), the coefficients of absorption of induced emission and of spontaneous emission are all equal. This is the content of the relation between the Einstein coefficients A and B of Chapter 42, Vol. I, Eq. (42.18).

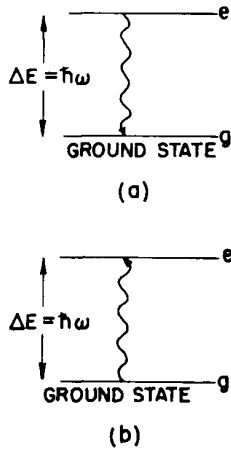


Fig. 4-6. Radiation and absorption of a photon with the frequency ω .

4-5 The blackbody spectrum

We would like to use our rules for Bose particles to discuss once more the spectrum of blackbody radiation (see Chapter 42, Vol. I). We will do it by finding out how many photons there are in a box if the radiation is in thermal equilibrium with some atoms in the box. Suppose that for each light frequency ω , there are a certain number N of atoms which have two energy states separated by the energy $\Delta E = \hbar\omega$. See Fig. 4-6. We'll call the lower-energy state the "ground" state and the upper state the "excited" state. Let N_g and N_e be the average numbers of atoms in the ground and excited states; then in thermal equilibrium at the temperature T , we have from statistical mechanics that

$$\frac{N_e}{N_g} = e^{-\Delta E/kT} = e^{-\hbar\omega/kT}. \quad (4.30)$$

Each atom in the ground state can absorb a photon and go into the excited state, and each atom in the excited state can emit a photon and go to the ground state. In equilibrium, the rates for these two processes must be equal. The rates are proportional to the probability for the event and to the number of atoms present. Let's let \bar{n} be the average number of photons present in a given state with the frequency ω . Then the absorption rate from that state is $N_g \bar{n} |a|^2$, and the emission rate into that state is $N_e (\bar{n} + 1) |a|^2$. Setting the two rates equal, we have that

$$N_g \bar{n} = N_e (\bar{n} + 1). \quad (4.31)$$

Combining this with Eq. (4.30), we have

$$\frac{\bar{n}}{\bar{n} + 1} = e^{-\hbar\omega/kT}.$$

Solving for \bar{n} , we have

$$\bar{n} = \frac{1}{e^{\hbar\omega/kT} - 1}, \quad (4.32)$$

which is the mean number of photons in any state with frequency ω , for a cavity in thermal equilibrium. Since each photon has the energy $\hbar\omega$, the energy in the photons of a given state is $\bar{n}\hbar\omega$, or

$$\frac{\hbar\omega}{e^{\hbar\omega/kT} - 1}. \quad (4.33)$$

Incidentally, we once found a similar equation in another context [Chapter 41, Vol. I, Eq. (41.15)]. You remember that for any harmonic oscillator—such as a weight on a spring—the quantum mechanical energy levels are equally spaced with a separation $\hbar\omega$, as drawn in Fig. 4-7. If we call the energy of the n th level $n\hbar\omega$, we find that the mean energy of such an oscillator is also given by Eq. (4.33). Yet this equation was derived here for photons, by counting particles, and it gives the same results. That is one of the marvelous miracles of quantum mechanics. If one begins by considering a kind of state or condition for Bose particles which do not interact with each other (we have assumed that the photons do not interact with each other), and then considers that into this state there can be put either zero, or one, or two, . . . up to any number n of particles, one finds that this system behaves for all quantum mechanical purposes exactly like a harmonic oscillator. By such an oscillator we mean a dynamic system like a weight on a spring or a standing wave in a resonant cavity. And that is why it is possible to represent the electromagnetic field by photon particles. From one point of view, we can analyze the electromagnetic field in a box or cavity in terms of a lot of harmonic oscillators, treating each mode of oscillation according to quantum mechanics as a harmonic oscillator. From a different point of view, we can analyze the same physics in terms of identical Bose particles. And the results of both ways of working *are always in exact agreement*. There is no way to make up your mind whether the electromagnetic field is really to be described as a quantized harmonic oscillator or by giving how many photons there are in each condition. The two views turn out to be mathematically identical. So in the future we can speak either about the number of photons in a particular state in a box or the number of the energy level associated with a particular mode of oscillation of the electromagnetic field. They are two ways of saying the same thing. The same is true of photons in free space. They are equivalent to oscillations of a cavity whose walls have receded to infinity.

We have computed the mean energy in any particular mode in a box at the temperature T ; we need only one more thing to get the blackbody radiation law: We need to know how many modes there are at each energy. (We assume that for every mode there are some atoms in the box—or in the walls—which have energy levels that can radiate into that mode, so that each mode can get into thermal equilibrium.) The blackbody radiation law is usually stated by giving the energy per unit volume carried by the light in a small frequency interval from ω to $\omega + \Delta\omega$. So we need to know how many modes there are in a box with frequencies in the

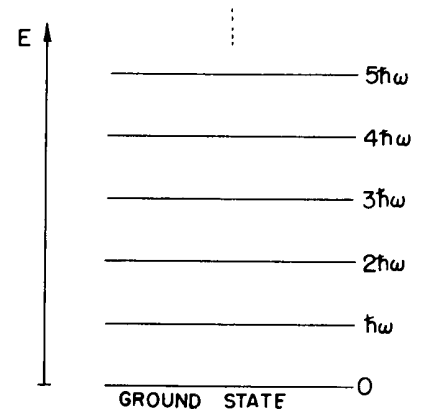


Fig. 4-7. The energy levels of a harmonic oscillator.

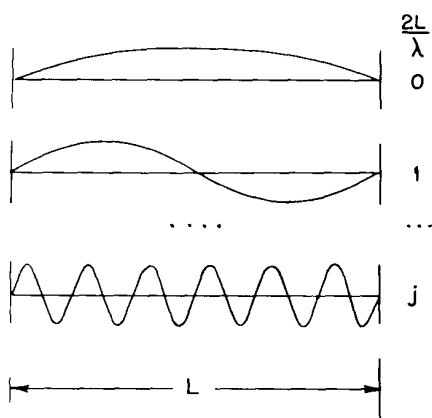


Fig. 4-8. The standing wave modes on a line.

interval $\Delta\omega$. Although this question continually comes up in quantum mechanics, it is purely a classical question about standing waves.

We will get the answer only for a rectangular box. It comes out the same for a box of any shape, but it's very complicated to compute for the arbitrary case. Also, we are only interested in a box whose dimensions are very large compared with a wavelength of the light. Then there are billions and billions of modes; there will be many in any small frequency interval $\Delta\omega$, so we can speak of the "average number" in any $\Delta\omega$ at the frequency ω . Let's start by asking how many modes there are in a one-dimensional case—as for waves on a stretched string. You know that each mode is a sine wave that has to go to zero at both ends; in other words, there must be an integral number of half-wavelengths in the length of the line, as shown in Fig. 4-8. We prefer to use the wave number $k = 2\pi/\lambda$; calling k_j the wave number of the j th mode, we have that

$$k_j = \frac{j\pi}{L}, \quad (4.34)$$

where j is any integer. The separation δk between successive modes is

$$\delta k = k_{j+1} - k_j = \frac{\pi}{L}.$$

We want to assume that kL is so large that in a small interval Δk , there are many modes. Calling $\Delta\mathfrak{N}$ the number of modes in the interval Δk , we have

$$\Delta\mathfrak{N} = \frac{\Delta k}{\delta k} = \frac{L}{\pi} \Delta k. \quad (4.35)$$

Now theoretical physicists working in quantum mechanics usually prefer to say that there are one-half as many modes; they write

$$\Delta\mathfrak{N} = \frac{L}{2\pi} \Delta k. \quad (4.36)$$

We would like to explain why. They usually like to think in terms of travelling waves—some going to the right (with a positive k) and some going to the left (with a negative k). But a "mode" is a *standing* wave which is the sum of two waves, one going in each direction. In other words, they consider each standing wave as containing two distinct photon "states." So if by $\Delta\mathfrak{N}$, one prefers to mean the number of photon states of a given k (where now k ranges over positive and negative values), one should then take $\Delta\mathfrak{N}$ half as big. (All integrals must now go from $k = -\infty$ to $k = +\infty$, and the total number of states up to any given absolute value of k will come out O.K.) Of course, we are not then describing standing waves very well, but we are counting modes in a consistent way.

Now we want to extend the results to three dimensions. A standing wave in a rectangular box must have an integral number of half-waves *along each axis*. The situation for two of the dimensions is shown in Fig. 4-9. Each wave direction and frequency is described by a vector wave number \mathbf{k} , whose x , y , and z components must satisfy equations like Eq. (4.34). So we have that

$$k_x = \frac{j_x\pi}{L_x},$$

$$k_y = \frac{j_y\pi}{L_y},$$

$$k_z = \frac{j_z\pi}{L_z}.$$

The number of modes with k_x in an interval Δk_x is, as before,

$$\frac{L_x}{2\pi} \Delta k_x,$$

and similarly for Δk_y and Δk_z . If we call $\Delta\mathfrak{N}(\mathbf{k})$ the number of modes for a vector

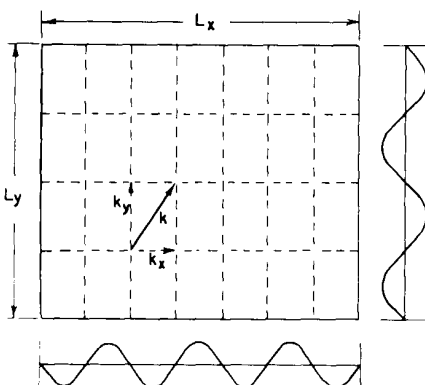


Fig. 4-9. Standing wave modes in two dimensions.

wave number \mathbf{k} whose x -component is between k_x and $k_x + \Delta k_x$, whose y -component is between k_y and $k_y + \Delta k_y$, and whose z -component is between k_z and $k_z + \Delta k_z$, then

$$\Delta \mathcal{N}(\mathbf{k}) = \frac{L_x L_y L_z}{(2\pi)^3} \Delta k_x \Delta k_y \Delta k_z. \quad (4.37)$$

The product $L_x L_y L_z$ is equal to the volume V of the box. So we have the important result that for high frequencies (wavelengths small compared with the dimensions), the number of modes in a cavity is proportional to the volume V of the box and to the “volume in k -space” $\Delta k_x \Delta k_y \Delta k_z$. This result comes up again and again in many problems and should be memorized:

$$d\mathcal{N}(\mathbf{k}) = V \frac{d^3 \mathbf{k}}{(2\pi)^3}. \quad (4.38)$$

Although we have not proved it, the result is independent of the shape of the box.

We will now apply this result to find the number of photon modes for photons with frequencies in the range $\Delta\omega$. We are just interested in the energy in various modes—but not interested in the directions of the waves. We would like to know the number of modes in a given range of frequencies. In a vacuum the magnitude of \mathbf{k} is related to the frequency by

$$|\mathbf{k}| = \frac{\omega}{c}. \quad (4.39)$$

So in a frequency interval $\Delta\omega$, these are all the modes which correspond to \mathbf{k} 's with a *magnitude* between k and $k + \Delta k$, independent of the direction. The “volume in k -space” between k and $k + \Delta k$ is a spherical shell of volume

$$4\pi k^2 \Delta k.$$

The number of modes is then

$$\Delta \mathcal{N}(\omega) = \frac{V 4\pi k^2 \Delta k}{(2\pi)^3}. \quad (4.40)$$

However, since we are now interested in frequencies, we should substitute $k = \omega/c$, so we get

$$\Delta \mathcal{N}(\omega) = \frac{V 4\pi \omega^2 \Delta\omega}{(2\pi)^3 c^3}. \quad (4.41)$$

There is one more complication. If we are talking about modes of an electromagnetic wave, for any given wave vector \mathbf{k} there can be either of two polarizations (at right angles to each other). Since these modes are independent, we must—*for light*—double the number of modes. So we have

$$\Delta \mathcal{N}(\omega) = \frac{V \omega^2 \Delta\omega}{\pi^2 c^3} \quad (\text{for light}). \quad (4.42)$$

We have shown, Eq. (4.33), that each mode (or each “state”) has on the average the energy

$$\bar{n} \hbar \omega = \frac{\hbar \omega}{e^{\hbar \omega / kT} - 1}.$$

Multiplying this by the number of modes, we get the energy ΔE in the modes that lie in the interval $\Delta\omega$:

$$\Delta E = \frac{\hbar \omega}{e^{\hbar \omega / kT} - 1} \frac{V \omega^2 \Delta\omega}{\pi^2 c^3}. \quad (4.43)$$

This is the law for the frequency spectrum of blackbody radiation, which we have already found in Chapter 41 of Vol. I. The spectrum is plotted in Fig. 4-10. You see now that the answer depends on the fact that photons are Bose particles, which

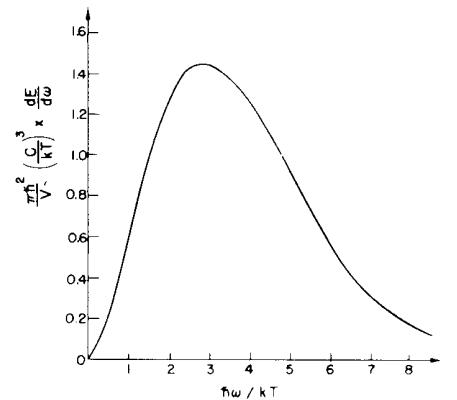


Fig. 4-10. The frequency spectrum of radiation in a cavity in thermal equilibrium, the “blackbody” spectrum.

have a tendency to try to get all into the same state (because the amplitude for doing so is large). You will remember, it was Planck's study of the blackbody spectrum (which was a mystery to classical physics), and his discovery of the formula in Eq. (4.43) that started the whole subject of quantum mechanics.

4-6 Liquid helium

Liquid helium has at low temperatures many odd properties which we cannot unfortunately take the time to describe in detail right now, but many of them arise from the fact that a helium atom is a Bose particle. One of the things is that liquid helium flows without any viscous resistance. It is, in fact, the ideal "dry" water we have been talking about in one of the earlier chapters—provided that the velocities are low enough. The reason is the following. In order for a liquid to have viscosity, there must be internal energy losses; there must be some way for one part of the liquid to have a motion that is different from that of the rest of the liquid. This means that it must be possible to knock some of the atoms into states that are different from the states occupied by other atoms. But at sufficiently low temperatures, when the thermal motions get very small, all the atoms try to get into the same condition. So, if some of them are moving along, then all the atoms try to move together in the same state. There is a kind of rigidity to the motion, and it is hard to break the motion up into irregular patterns of turbulence, as would happen, for example, with independent particles. So in a liquid of Bose particles, there is a strong tendency for all the atoms to go into the same state—which is represented by the $\sqrt{n+1}$ factor we found earlier. (For a bottle of liquid helium n is, of course, a very large number!) This cooperative motion does not happen at high temperatures, because then there is sufficient thermal energy to put the various atoms into various different higher states. But at a sufficiently low temperature there suddenly comes a moment in which all the helium atoms try to go into the same state. The helium becomes a superfluid. Incidentally, this phenomenon only appears with the isotope of helium which has atomic weight 4. For the helium isotope of atomic weight 3, the individual atoms are Fermi particles, and the liquid is a normal fluid. Since superfluidity occurs only with He^4 , it is evidently a quantum mechanical effect—due to the Bose nature of the α -particle.

4-7 The exclusion principle

Fermi particles act in a completely different way. Let's see what happens if we try to put two Fermi particles into the same state. We will go back to our original example and ask for the amplitude that two identical Fermi particles will be scattered into almost exactly the same direction. The amplitude that particle a will go in direction 1 and particle b will go in direction 2 is

$$\langle 1 | a \rangle \langle 2 | b \rangle,$$

whereas the amplitude that the outgoing directions will be interchanged is

$$\langle 2 | a \rangle \langle 1 | b \rangle.$$

Since we have Fermi particles, the amplitude for the process is the difference of these two amplitudes:

$$\langle 1 | a \rangle \langle 2 | b \rangle - \langle 2 | a \rangle \langle 1 | b \rangle. \quad (4.44)$$

Let's say that by "direction 1" we mean that the particle has not only a certain direction but also a given direction of its spin, and that "direction 2" is almost exactly the same as direction 1 and corresponds to the *same* spin direction. Then $\langle 1 | a \rangle$ and $\langle 2 | a \rangle$ are nearly equal. (This would not necessarily be true if the outgoing states 1 and 2 did not have the same spin, because there might be some reason why the amplitude would depend on the spin direction.) Now if we let

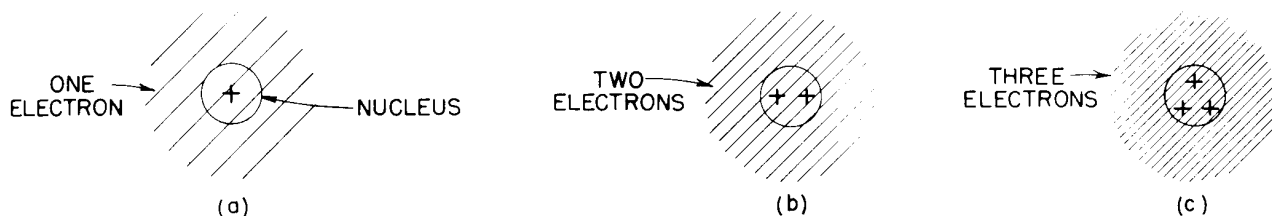


Fig. 4-11. How atoms might look if electrons behaved like Bose particles.

directions 1 and 2 approach each other, the total amplitude in Eq. (4.44) becomes zero. The result for Fermi particles is much simpler than for Bose particles. It just isn't possible at all for two Fermi particles—such as two electrons—to get into exactly the same state. You will never find two electrons in the same position with their two spins in the same direction. It is not possible for two electrons to have the same momentum and the same spin directions. If they are at the same location or with the same state of motion, the only possibility is that they must be spinning opposite to each other.

What are the consequences of this? There are a number of most remarkable effects which are a consequence of the fact that two Fermi particles cannot get into the same state. In fact, almost all the peculiarities of the material world hinge on this wonderful fact. The variety that is represented in the periodic table is basically a consequence of this one rule.

Of course, we cannot say what the world would be like if this one rule were changed, because it is just a part of the whole structure of quantum mechanics, and it is impossible to say what else would change if the rule about Fermi particles were different. Anyway, let's just try to see what would happen if only this one rule were changed. First, we can show that every atom would be more or less the same. Let's start with the hydrogen atom. It would not be noticeably affected. The proton of the nucleus would be surrounded by a spherically symmetric electron cloud, as shown in Fig. 4-11(a). As we have described in Chapter 2, the electron is attracted to the center, but the uncertainty principle requires that there be a balance between the concentration in space and in momentum. The balance means that there must be a certain energy and a certain spread in the electron distribution which determines the characteristic dimension of the hydrogen atom.

Now suppose that we have a nucleus with two units of charge, such as the helium nucleus. This nucleus would attract two electrons, and if they were Bose particles, they would—except for their electric repulsion—both crowd in as close as possible to the nucleus. A helium atom might look as shown in part (b) of the figure. Similarly, a lithium atom which has a triply charged nucleus would have an electron distribution like that shown in part (c) of Fig. 4-11. Every atom would look more or less the same—a little round ball with all the electrons sitting near the nucleus, nothing directional and nothing complicated.

Because electrons are Fermi particles, however, the actual situation is quite different. For the hydrogen atom the situation is essentially unchanged. The only difference is that the electron has a spin which we indicate by the little arrow in Fig. 4-12(a). In the case of a helium atom, however, we cannot put two electrons on top of each other. But wait, that is only true if their spins are the same. Two electrons *can* occupy the same state if their spins are opposite. So the helium atom does not look much different either. It would appear as shown in part (b) of Fig. 4-12. For lithium, however, the situation becomes quite different. Where can we put the third electron? The third electron cannot go on top of the other two because both spin directions are occupied. (You remember that for an electron or any particle with spin $1/2$ there are only two possible directions for the spin.) The third electron can't go near the place occupied by the other two, so it must take up a special condition in a different kind of state farther away from the nucleus in part (c) of the figure. (We are speaking only in a rather rough way here, because in reality all three electrons are identical; since we cannot really distinguish which one is which, our picture is only an approximate one.)

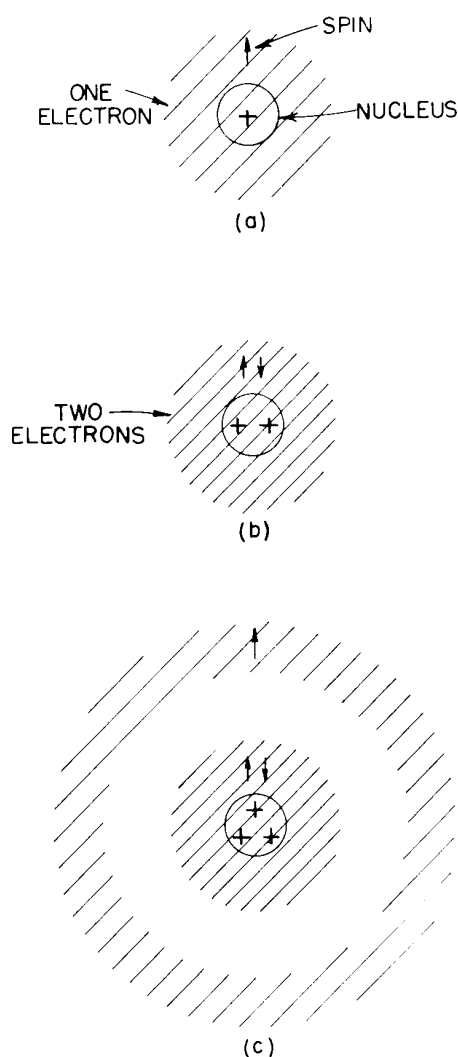


Fig. 4-12. Atomic configurations for real, Fermi-type, spin one-half electrons.

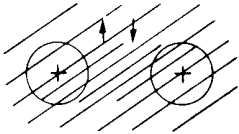


Fig. 4-13. The hydrogen molecule.

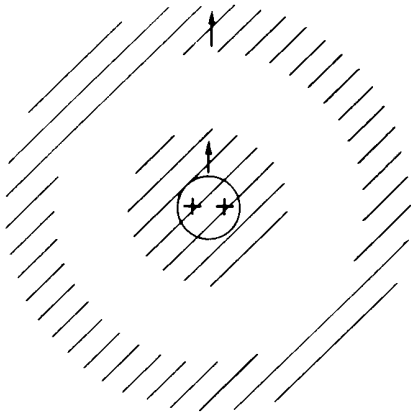


Fig. 4-14. Helium with one electron in a higher energy state.

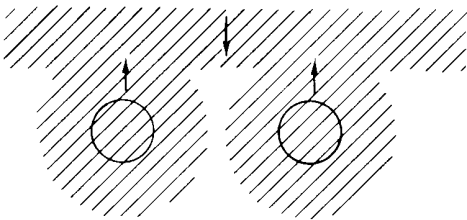


Fig. 4-15. The likely mechanism in a ferromagnetic crystal; the conduction electron is antiparallel to the unpaired inner electrons.

Now we can begin to see why different atoms will have different chemical properties. Because the third electron in lithium is farther out, it is relatively more loosely bound. It is much easier to remove an electron from lithium than from helium. (Experimentally, it takes 25 volts to ionize helium but only 5 volts to ionize lithium.) This accounts for the valence of the lithium atom. The directional properties of the valence have to do with the pattern of the waves of the outer electron, which we will not go into at the moment. But we can already see the importance of the so-called *exclusion principle*—which states that no two electrons can be found in exactly the same state (including spin).

The exclusion principle is also responsible for the stability of matter on a large scale. We explained earlier that the individual atoms in matter did not collapse because of the uncertainty principle; but this does not explain why it is that two hydrogen atoms can't be squeezed together as close as you want—why it is that all the protons don't get close together with one big smear of electrons around them. The answer is, of course, that since no more than two electrons—with opposite spins—can be in roughly the same place, the hydrogen atoms must keep away from each other. So the stability of matter on a large scale is really a consequence of the Fermi particle nature of the electrons.

Of course, if the outer electrons on two atoms have spins in opposite directions, they can get close to each other. This is, in fact, just the way that the chemical bond comes about. It turns out that two atoms together will generally have the lowest energy if there is an electron between them. It is a kind of an electrical attraction for the two positive nuclei toward the electron in the middle. It is possible to put two electrons more or less between the two nuclei so long as their spins are opposite, and the strongest chemical binding comes about this way. There is no stronger binding, because the exclusion principle does not allow there to be more than two electrons in the space between the atoms. We expect the hydrogen molecule to look more or less as shown in Fig. 4-13.

We want to mention one more consequence of the exclusion principle. You remember that if both electrons in the helium atom are to be close to the nucleus, their spins are necessarily opposite. Now suppose that we would like to try to arrange to have both electrons with the same spin—as we might consider doing by putting on a fantastically strong magnetic field that would try to line up the spins in the same direction. But then the two electrons could not occupy the same state in space. One of them would have to take on a different geometrical position, as indicated in Fig. 4-14. The electron which is located farther from the nucleus has less binding energy. The energy of the whole atom is therefore quite a bit higher. In other words, when the two spins are opposite, there is a much stronger total attraction.

So, there is an apparent, enormous force trying to line up spins opposite to each other when two electrons are close together. If two electrons are trying to go in the same place, there is a very strong tendency for the spins to become lined opposite. This apparent force trying to orient the two spins opposite to each other is much more powerful than the tiny force between the two magnetic moments of the electrons. You remember when we were speaking of ferromagnetism there was the mystery of why the electrons in different atoms had a strong tendency to line up parallel. Although there is still no quantitative explanation, it is believed that what happens is that the electrons around the core of one atom interact through the exclusion principle with the outer electrons which have become free to wander throughout the crystal. This interaction causes the spins of the free electrons and the inner electrons to take on opposite directions. But the free electrons and the inner atomic electrons can only be opposite provided all the inner electrons have the same spin direction, as indicated in Fig. 4-15. It seems probable that it is the effect of the exclusion principle acting indirectly through the free electrons that gives rise to the strong aligning forces responsible for ferromagnetism.

We will mention one further example of the influence of the exclusion principle. We have said earlier that the nuclear forces are the same between the neutron and the proton, between the proton and the proton, and between the proton and the neutron. Why is it then that a proton and a neutron can stick together to make a

deuterium nucleus, whereas there is no nucleus with just two protons or with just two neutrons? The deuteron is, as a matter of fact, bound by an energy of about 2.2 million volts, yet, there is no corresponding binding between a pair of protons to make an isotope of helium with the atomic weight 2. Such nuclei do not exist. The combination of two protons does not make a bound state.

The answer is a result of two effects: first, the exclusion principle; and second, the fact that the nuclear forces are somewhat sensitive to the direction of spin. The force between a neutron and a proton is attractive and somewhat stronger when the spins are parallel than when they are opposite. It happens that these forces are just different enough that a deuteron can only be made if the neutron and proton have their spins parallel; when their spins are opposite, the attraction is not quite strong enough to bind them together. Since the spins of the neutron and proton are each one-half and are in the same direction, the deuteron has a spin of one. We know, however, that two protons are not allowed to sit on top of each other if their spins are parallel. If it were not for the exclusion principle, two protons would be bound, but since they cannot exist at the same place and with the same spin directions, the He^2 nucleus does not exist. The protons could come together with their spins opposite, but then there is not enough binding to make a stable nucleus, because the nuclear force for opposite spins is too weak to bind a pair of nucleons. The attractive force between neutrons and protons of opposite spins can be seen by scattering experiments. Similar scattering experiments with two protons with parallel spins show that there is the corresponding attraction. So it is the exclusion principle that helps explain why deuterium can exist when He^2 cannot.

Spin One

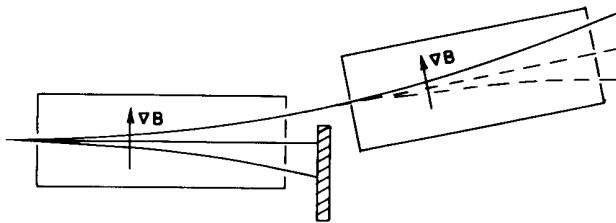
5-1 Filtering atoms with a Stern-Gerlach apparatus

In this chapter we really begin the quantum mechanics proper—in the sense that we are going to describe a quantum mechanical phenomenon in a completely quantum mechanical way. We will make no apologies and no attempt to find connections to classical mechanics. We want to talk about something new in a new language. The particular situation which we are going to describe is the behavior of the so-called quantization of the angular momentum, for a particle of *spin one*. But we won't use words like "angular momentum" or other concepts of classical mechanics until later. We have chosen this particular example because it is relatively simple, although not the simplest possible example. It is sufficiently complicated that it can stand as a prototype which can be generalized for the description of all quantum mechanical phenomena. Thus, although we are dealing with a particular example, all the laws which we mention are immediately generalizable, and we will give the generalizations so that you will see the general characteristics of a quantum mechanical description. We begin with the phenomenon of the splitting of a beam of atoms into three separate beams in a Stern-Gerlach experiment.

You remember that if we have an inhomogeneous magnetic field made by a magnet with a pointed pole tip and we send a beam through the apparatus, the beam of particles may be split into a number of beams—the number depending on the particular kind of atom and its state. We are going to take the case of an atom which gives three beams, and we are going to call that a particle of *spin one*. You can do for yourself the case of five beams, seven beams, two beams, etc.—you just copy everything down and where we have three terms, you will have five terms, seven terms, and so on.

Imagine the apparatus drawn schematically in Fig. 5-1. A beam of atoms (or particles of any kind) is collimated by some slits and passes through a non-uniform field. Let's say that the beam moves in the y -direction and that the magnetic field and its gradient are both in the z -direction. Then, looking from the side, we will see the beam split vertically into three beams, as shown in the figure. Now at the output end of the magnet we could put small counters which count the rate of arrival of particles in any one of the three beams. Or we can block off two of the beams and let the third one go on.

Suppose we block off the lower two beams and let the top-most beam go on and enter a second Stern-Gerlach apparatus of the same kind, as shown in Fig. 5-2. What happens? There are *not* three beams in the second apparatus; there is only the top beam.† This is what you would expect if you think of the second apparatus as simply an extension of the first. Those atoms which are being pushed upward continue to be pushed upward in the second magnet.



† We are assuming that the deflection angles are very small.

5-1 Filtering atoms with a Stern-Gerlach apparatus

5-2 Experiments with filtered atoms

5-3 Stern-Gerlach filters in series

5-4 Base states

5-5 Interfering amplitudes

5-6 The machinery of quantum mechanics

5-7 Transforming to a different base

5-8 Other situations

Review: Chapter 35, Vol. II, *Paramagnetism and Magnetic Resonance*. For your convenience this chapter is reproduced in the Appendix of this volume.

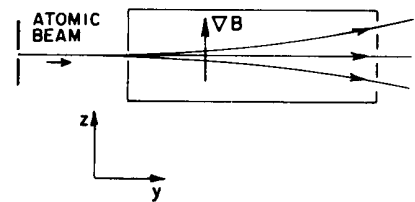


Fig. 5-1. In a Stern-Gerlach experiment, atoms of spin one are split into three beams.

Fig. 5-2. The atoms from one of the beams are sent into a second identical apparatus.

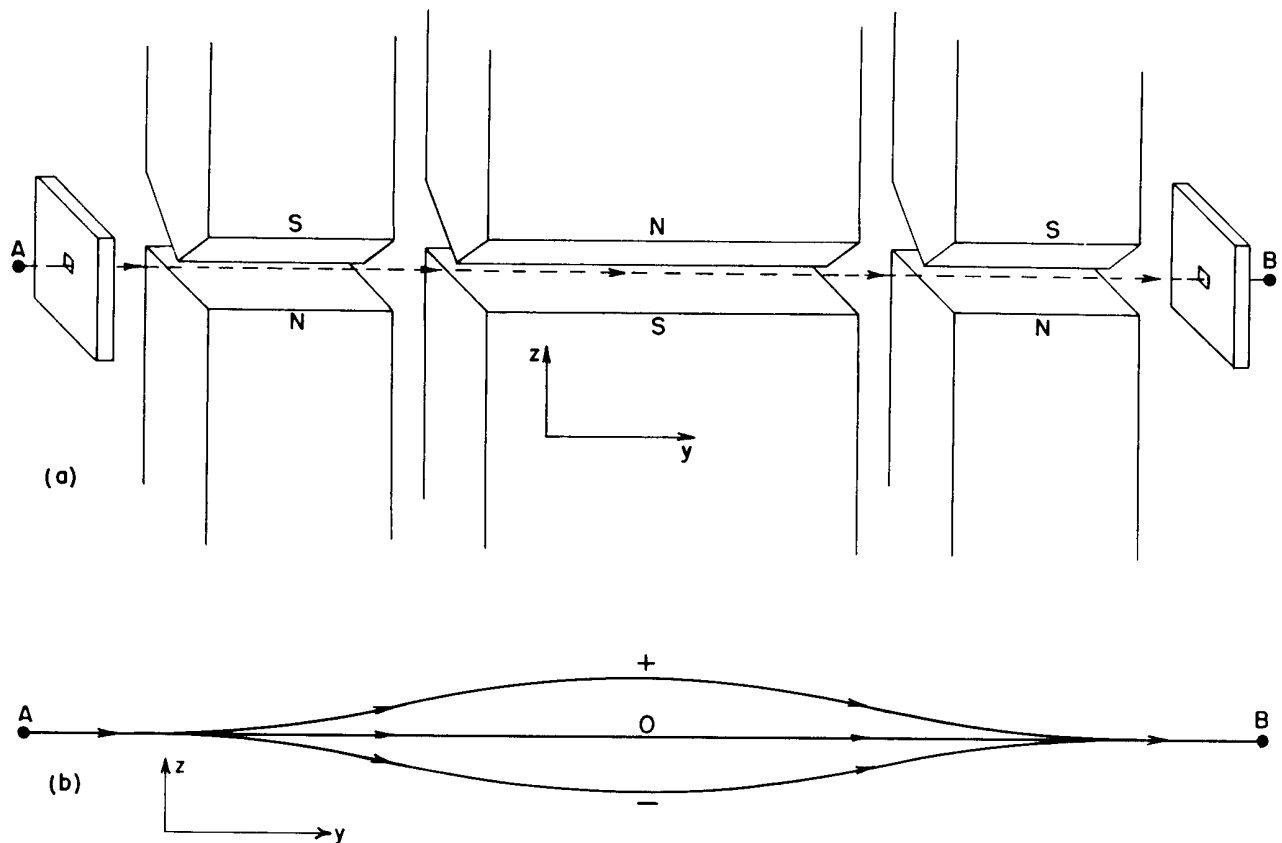


Fig. 5-3. (a) An imagined modification of a Stern-Gerlach apparatus. (b) The paths of spin-one atoms.

You can see then that the first apparatus has produced a beam of “purified” objects—atoms that get bent upward in the particular inhomogeneous field. The atoms, as they enter the original Stern-Gerlach apparatus, are of three “varieties,” and the three kinds take different trajectories. By filtering out all but one of the varieties, we can make a beam whose future behavior in the same kind of apparatus is determined and predictable. We will call this a *filtered* beam, or a *polarized* beam, or a beam in which the atoms all are known to be in a *definite state*.

For the rest of our discussion, it will be more convenient if we consider a somewhat modified apparatus of the Stern-Gerlach type. The apparatus looks more complicated at first, but it will make all the arguments simpler. Anyway, since they are only “thought experiments,” it doesn’t cost anything to complicate the equipment. (Incidentally, no one has ever done all of the experiments we will describe in just this way, but we know what *would* happen from the laws of quantum mechanics, which are, of course, based on other similar experiments. These other experiments are harder to understand at the beginning, so we want to describe some idealized—but possible—experiments.)

Figure 5-3(a) shows a drawing of the “modified Stern-Gerlach apparatus” we would like to use. It consists of a sequence of three high-gradient magnets. The first one (on the left) is just the usual Stern-Gerlach magnet and splits the incoming beam of spin-one particles into three separate beams. The second magnet has the same cross section as the first, but is twice as long *and* the polarity of its magnetic field is opposite the field in magnet 1. The second magnet pushes in the opposite direction on the atomic magnets and bends their paths back toward the axis, as shown in the trajectories drawn in the lower part of the figure. The third magnet is just like the first, and brings the three beams back together again, so that leaves the exit hole along the axis. Finally, we would like to imagine that in front of the hole at *A* there is some mechanism which can get the atoms started from rest and that after the exit hole at *B* there is a decelerating mechanism that brings the atoms back to rest at *B*. That is not essential, but it will mean that in

our analysis we won't have to worry about including any effects of the motion as the atoms come out, and can concentrate on those matters having only to do with the spin. The whole purpose of the "improved" apparatus is just to bring all the particles to the same place, and with zero speed.

Now if we want to do an experiment like the one in Fig. 5-2, we can first make a filtered beam by putting a plate in the middle of the apparatus that blocks two of the beams, as shown in Fig. 5-4. If we now put the polarized atoms through a second identical apparatus, all the atoms will take the upper path, as can be verified by putting similar plates in the way of the various beams of the second S filter and seeing whether particles get through.

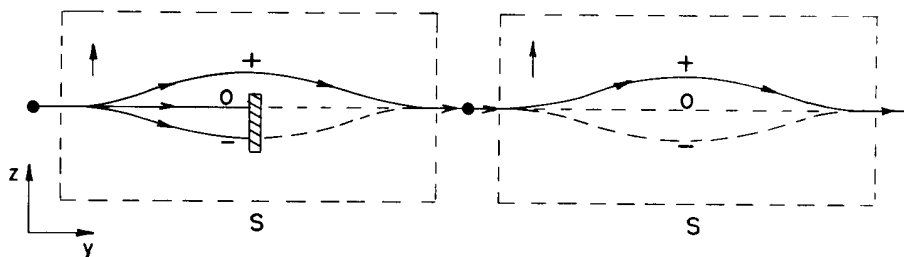


Fig. 5-4. The "improved" Stern-Gerlach apparatus as a filter.

Suppose we call the first apparatus by the name S . (We are going to consider all sorts of combinations, and we will need labels to keep things straight.) We will say that the atoms which take the top path in S are in the "plus state with respect to S "; the ones which take the middle path are in the "zero state with respect to S "; and the ones which take the lowest path are in the "minus state with respect to S ." (In the more usual language we would say that the z -component of the angular momentum was $+\hbar$, 0 , and $-\hbar$, but we are not using that language now.) Now in Fig. 5-4 the second apparatus is oriented just like the first, so the filtered atoms will all go on the upper path. Or if we had blocked off the upper and lower beams in the first apparatus and let only the zero state through, all the filtered atoms would go through the middle path of the second apparatus. And if we had blocked off all but the lowest beam in the first, there would be only a low beam in the second. We can say that in each case our first apparatus has produced a filtered beam in a *pure* state with respect to S (+, 0, or -), and we can test which state is present by putting the atoms through a second, identical apparatus.

We can make our second apparatus so that it transmits only atoms of a particular state—by putting masks inside it as we did for the first one—and then we can test the state of the incoming beam just by seeing whether anything comes out the far end. For instance, if we block off the two lower paths in the second apparatus, 100 percent of the atoms will still come through; but if we block off the upper path, nothing will get through.

To make this kind of discussion easier, we are going to invent a shorthand symbol to represent one of our improved Stern-Gerlach apparatuses. We will let the symbol

$$\left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_S \quad (5.1)$$

stand for one complete apparatus. (This is *not* a symbol you will ever find used in quantum mechanics; we've just invented it for this chapter. It is simply meant to be a shorthand picture of the apparatus of Fig. 5-3.) Since we are going to want to use several apparatuses at once, and with various orientations, we will identify each with a letter underneath. So the symbol in (5.1) stands for the apparatus S . When we block off one or more of the beams inside, we will show that by some

vertical bars indicating which beam is blocked, like this:

$$\left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_s. \quad (5.2)$$


The various possible combinations we will be using are shown in Fig. 5-5.

If we have two filters in succession (as in Fig. 5-4), we will put the two symbols next to each other, like this:

$$\left\{ \begin{matrix} + \\ 0 \\ - \end{matrix} \right\}_{\mathcal{S}} = \left\{ \begin{matrix} + \\ 0 \\ - \end{matrix} \right\}_{\mathcal{S}}. \quad (5.3)$$

For this setup, everything that comes through the first also gets through the second. In fact, even if we block off the “zero” and “minus” channels of the second apparatus, so that we have

$$\left\{ \begin{array}{c} + \\ o \\ - \end{array} \right\} = \text{(a) } \left[\begin{array}{c} \uparrow \\ \text{---} \end{array} \right] \quad \left\{ \begin{array}{c} + \\ o \\ - \end{array} \right\}_S \quad \left\{ \begin{array}{c} + \\ o \\ - \end{array} \right\}_S, \quad (5.4)$$

(b)  we still get 100 percent transmission through the second apparatus. On the other hand, if we have

$$\left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_S = \left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_S, \quad (5.5)$$

nothing at all comes out of the far end. Similarly,


(d) 

Fig. 5-5. Special shorthand symbols for Stern-Gerlach type filters.

Fig. 5-5. Special shorthand symbols for Stern-Gerlach type filters.

would be just equivalent to

$$\left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}$$

by itself.

Now we want to describe these experiments quantum mechanically. We will say that an atom is in the $(+S)$ state if it has gone through the apparatus of Fig. 5-5(b), that it is in a $(0S)$ state if it has gone through (c), and in a $(-S)$ state if it has gone through (d).† Then we let $\langle b | a \rangle$ be the *amplitude* that an atom which is in state a will get through an apparatus into the b state. We can say: $\langle b | a \rangle$ is the amplitude for an atom *in* the state a to *get into* the state b . The experiment (5.4) gives us that

$$\langle +S | +S \rangle = 1,$$

† Read: $(+S)$ = “plus- S ”; $(0S)$ = “zero- S ”; $(-S)$ = “minus- S .”

whereas (5.5) gives us

$$\langle -S | +S \rangle = 0.$$

Similarly, the result of (5.6) is

$$\langle +S | -S \rangle = 0,$$

and of (5.7) is

$$\langle -S | -S \rangle = 1.$$

As long as we deal only with “pure” states—that is, we have only one channel open—there are nine such amplitudes, and we can write them in a table:

$$\begin{array}{c} \text{to} \\ \begin{array}{c} +S \\ 0S \\ -S \end{array} \end{array} \begin{array}{c} \text{from} \\ \begin{array}{ccc} +S & 0S & -S \end{array} \\ \hline \begin{array}{ccc} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{array} \end{array} \quad (5.8)$$

This array of nine numbers—called a *matrix*—summarizes the phenomena we’ve been describing.

5-2 Experiments with filtered atoms

Now comes the big question: What happens if the second apparatus is tipped to a different angle, so that its field axis is no longer parallel to the first? It could be not only tipped, but also pointed in a different direction—for instance, it could take the beam off at 90° with respect to the original direction. To take it easy at first, let’s first think about an arrangement in which the second Stern-Gerlach experiment is tilted by some angle α about the y -axis, as shown in Fig. 5-6. We’ll call the second apparatus T . Suppose that we now set up the following experiment:

$$\begin{array}{c} \left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\} \\ s \end{array} \quad \begin{array}{c} \left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\} \\ t \end{array},$$

or the experiment:

$$\begin{array}{c} \left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\} \\ s \end{array} \quad \begin{array}{c} \left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\} \\ t \end{array}.$$

What comes out at the far end in these cases?

The answer is this: If the atoms are in a definite state with respect to S , they are *not* in the same state with respect to T —a $(+S)$ state is *not* also a $(+T)$ state. There *is*, however, a certain *amplitude* to find the atom in a $(+T)$ state—or a $(0T)$ state or a $(-T)$ state.

In other words, as careful as we have been to make sure that we have the atoms in a definite condition, the fact of the matter is that if it goes through an apparatus which is tilted at a different angle it has, so to speak, to “reorient”

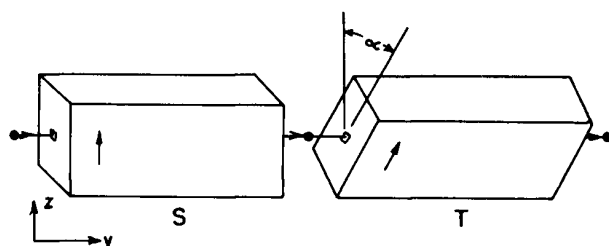


Fig. 5-6. Two Stern-Gerlach type filters in series; the second is tilted at the angle α with respect to the first.

itself—which it does, don't forget, by luck. We can put only one particle through at a time, and then we can only ask the question: What is the probability that it gets through? Some of the atoms that have gone through S will end in a $(+T)$ state, some of them will end in a $(0T)$, and some in a $(-T)$ state—all with different odds. These odds can be calculated by the absolute squares of complex amplitudes; what we want is some mathematical method, or quantum mechanical description, for these amplitudes. What we need to know are various quantities like

$$\langle -T | +S \rangle,$$

by which we mean the amplitude that an atom initially in the $(+S)$ state can get into the $(-T)$ condition (which is *not* zero unless T and S are lined up parallel to each other). There are other amplitudes like

$$\langle +T | 0S \rangle, \quad \text{or} \quad \langle 0T | -S \rangle, \quad \text{etc.}$$

There are, in fact, nine such amplitudes—another matrix—that a theory of particles should tell us how to calculate. Just as $F = ma$ tells us how to calculate what happens to a classical particle in any circumstance, the laws of quantum mechanics permit us to determine the amplitude that a particle will get through a particular apparatus. The central problem, then, is to be able to calculate—for any given tilt angle α , or in fact for any orientation whatever—the nine amplitudes:

$$\begin{array}{lll} \langle +T | +S \rangle, & \langle +T | 0S \rangle, & \langle +T | -S \rangle, \\ \langle 0T | +S \rangle, & \langle 0T | 0S \rangle, & \langle 0T | -S \rangle, \\ \langle -T | +S \rangle, & \langle -T | 0S \rangle, & \langle -T | -S \rangle. \end{array} \quad (5.9)$$

We can already figure out some relations among these amplitudes. First, according to our definitions, the absolute square

$$|\langle +T | +S \rangle|^2$$

is the *probability* that an atom in a $(+S)$ state will enter a $(+T)$ state. We will often find it more convenient to write such squares in the equivalent form

$$\langle +T | +S \rangle \langle +T | +S \rangle^*.$$

In the same notation the number

$$\langle 0T | +S \rangle \langle 0T | +S \rangle^*$$

is the probability that a particle in the $(+S)$ state will enter the $(0T)$ state, and

$$\langle -T | +S \rangle \langle -T | +S \rangle^*$$

is the probability that it will enter the $(-T)$ state. But the way our apparatuses are made, every atom which enters the T apparatus must be found in *some* one of the three states of the T apparatus—there's nowhere else for a given kind of atom to go. So the sum of the three probabilities we've just written must be equal to 100 percent. We have the relation

$$\begin{aligned} \langle +T | +S \rangle \langle +T | +S \rangle^* + \langle 0T | +S \rangle \langle 0T | +S \rangle^* \\ + \langle -T | +S \rangle \langle -T | +S \rangle^* = 1. \end{aligned} \quad (5.10)$$

There are, of course, two other such equations that we get if we start with a $(0S)$ or a $(-S)$ state. But they are all we can easily get, so we'll go on to some other general questions.

5-3 Stern-Gerlach filters in series

Here is an interesting question: Suppose we had atoms filtered into the $(+S)$ state, then we put them through a second filter, say into a $(0T)$ state, and *then* through *another* $+S$ filter. (We'll call the last filter S' just so we can distinguish

it from the first S -filter.) Do the atoms remember that they were once in a $(+S)$ state? In other words, we have the following experiment:

$$\left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_S \quad \left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_T \quad \left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_{S'}. \quad (5.11)$$

We want to know whether all those that get through T also get through S' . *They do not.* Once they have been filtered by T , they *do not remember* in any way that they were in a $(+S)$ state when they entered T . Note that the second S apparatus in (5.11) is oriented exactly the same as the first, so it is still an S -type filter. The states filtered by S' are, of course, still $(+S)$, $(0S)$, and $(-S)$.

The important point is this: *If the T filter passes only one beam*, the fraction that gets through the second S filter depends only on the setup of the T filter, and is completely independent of what precedes it. The fact that the same atoms were once sorted by an S filter has no influence whatever on what they will do once they have been sorted again into a pure beam by a T apparatus. From then on, the probability for getting into different states is the same no matter what happened before they got into the T apparatus.

As an example, let's compare the experiment of (5.11) with the following experiment:

$$\left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_S \quad \left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_T \quad \left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_{S'} \quad (5.12)$$

in which only the first S is changed. Let's say that the angle α (between S and T) is such that in experiment (5.11) one-third of the atoms that get through T also get through S' . In experiment (5.12), although there will, in general, be a different number of atoms coming through T , the *same fraction of these*—one-third—will also get through S' .

We can, in fact, show from what you have learned earlier that the fraction of the atoms that come out of T and get through any particular S' depends only on T and S' , not on anything that happened earlier. Let's compare experiment (5.12) with

$$\left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_S \quad \left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_T \quad \left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_{S'}. \quad (5.13)$$

The amplitude that an atom that comes out of S will also get through both T and S' is, for the experiments of (5.12),

$$\langle +S | 0T \rangle \langle 0T | 0S \rangle.$$

The corresponding probability is

$$|\langle +S | 0T \rangle \langle 0T | 0S \rangle|^2 = |\langle +S | 0T \rangle|^2 |\langle 0T | 0S \rangle|^2.$$

The probability for experiment (5.13) is

$$|\langle 0S | 0T \rangle \langle 0T | 0S \rangle|^2 = |\langle 0S | 0T \rangle|^2 |\langle 0T | 0S \rangle|^2.$$

The ratio is

$$\frac{|\langle 0S | 0T \rangle|^2}{|\langle +S | 0T \rangle|^2}$$

and depends only on T and S' , and not at all on which beam $(+S)$, $(0S)$, or $(-S)$ is selected by S . (The absolute numbers may go up and down together depending on how much gets through T .) We would, of course, find the same result if we compared the probabilities that the atoms would go into the plus or the minus

states with respect to S' , or the ratio of the probabilities to go into the zero or minus states.

In fact, since these ratios depend only on which beam is allowed to pass through T , and not on the selection made by the first S filter, it is clear that we would get the same result even if the last apparatus were not an S filter. If we use for the third apparatus—which we will now call R —one rotated by some arbitrary angle with respect to T , we would find that a ratio such as $|\langle 0 R | 0 T \rangle|^2 / |\langle + R | 0 T \rangle|^2$ was independent of which beam was passed by the first filter S .

5-4 Base states

These results illustrate one of the basic principles of quantum mechanics: Any atomic system can be separated by a filtering process into a certain set of what we will call *base states*, and the future behavior of the atoms in any single given base state depends only on the nature of the base state—it is independent of any previous history.† The base states depend, of course, on the filter used; for instance, the three states $(+T)$, $(0T)$, and $(-T)$ are one set of base states; the three states $(+S)$, $(0S)$, and $(-S)$ are another. There are any number of possibilities each as good as any other.

We should be careful to say that we are considering *good* filters which do indeed produce “pure” beams. If, for instance, our Stern-Gerlach apparatus didn’t produce a good separation of the three beams so that we could not separate them cleanly by our masks, then we could not make a complete separation into base states. We can tell if we have pure base states by seeing whether or not the beams can be split again in another filter of the same kind. If we have a pure $(+T)$ state, for instance, all the atoms will go through

$$\left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_T,$$

and none will go through

$$\left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_T,$$

or through

$$\left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_T.$$

Our statement about base states means that it is *possible* to filter to some pure state, so that no further filtering by an identical apparatus is possible.

We must also point out that what we are saying is exactly true only in rather idealized situations. In any real Stern-Gerlach apparatus, we would have to worry about diffraction by the slits that could cause some atoms to go into states corresponding to different angles, or about whether the beams might contain atoms with different excitations of their internal states, and so on. We have idealized the situation so that we are talking only about the states that are split in a magnetic field; we are ignoring things having to do with position, momentum, internal excitations, and the like. In general, one would need to consider also base states which are sorted out with respect to such things also. But to keep the concepts simple, we are considering only our set of three states, which is sufficient for the exact treatment of the idealized situation in which the atoms don’t get torn up in

† We do not intend the word “base state” to imply anything more than what is said here. They are not to be thought of as “basic” in any sense. We are using the word base with the thought of a *basis* for a description, somewhat in the sense that one speaks of “numbers to the base ten.”

going through the apparatus, or otherwise badly treated, and come to rest when they leave the apparatus.

You will note that we always begin our thought experiments by taking a filter with only one channel open, so that we start with some definite base state. We do this because atoms come out of a furnace in various states determined at random by the accidental happenings inside the furnace. (It gives what is called an “unpolarized” beam.) This randomness involves probabilities of the “classical” kind—as in coin tossing—which are different from the quantum mechanical probabilities we are worrying about now. Dealing with an unpolarized beam would get us into additional complications that are better to avoid until after we understand the behavior of polarized beams. So don’t try to consider at this point what happens if the *first* apparatus lets more than one beam through. (We will tell you how you can handle such cases at the end of the chapter.)

Let’s now go back and see what happens when we go from a base state for one filter to a base state for a different filter. Suppose we start again with

$$\left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_S \quad \left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_T.$$

The atoms which come out of T are in the base state $(0\ T)$ and have no memory that they were once in the state $(+S)$. Some people would say that in the filtering by T we have “lost the information” about the previous state $(+S)$ because we have “disturbed” the atoms when we separated them into three beams in the apparatus T . But that is not true. The past information is not lost by the *separation* into three beams, but by the *blocking masks* that are put in—as we can see by the following set of experiments.

We start with a $+S$ filter and will call N the number of atoms that come through it. If we follow this by a $0\ T$ filter, the number of atoms that come out is some fraction of the original number, say αN . If we then put another $+S$ filter, only some fraction β of these atoms will get to the far end. We can indicate this in the following way:

$$\left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_S \xrightarrow{N} \left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_T \xrightarrow{\alpha N} \left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_{S'} \xrightarrow{\beta \alpha N}. \quad (5.14)$$

If our third apparatus S' selected a different state, say the $(0\ S)$ state, a different fraction, say γ , would get through.† We would have

$$\left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_S \xrightarrow{N} \left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_T \xrightarrow{\alpha N} \left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_{S'} \xrightarrow{\gamma \alpha N}. \quad (5.15)$$

Now suppose we repeat these two experiments but remove all the masks from T . We would then find the remarkable results as follows:

$$\left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_S \xrightarrow{N} \left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_T \xrightarrow{N} \left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_{S'} \xrightarrow{N}, \quad (5.16)$$

$$\left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_S \xrightarrow{N} \left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_T \xrightarrow{N} \left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_{S'} \xrightarrow{0}. \quad (5.17)$$

† In terms of our earlier notation $\alpha = |\langle 0\ T | +\ S \rangle|^2$, $\beta = |\langle +\ S | 0\ T \rangle|^2$, and $\gamma = |\langle 0\ S | 0\ T \rangle|^2$.

All the atoms get through S' in the first case, but *none* in the second case! This is one of the great laws of quantum mechanics. That nature works this way is not self-evident, but the results we have given correspond for our idealized situation to the quantum mechanical behavior observed in innumerable experiments.

5-5 Interfering amplitudes

How can it be that in going from (5.15) to (5.17)—by *opening more channels*—we let *fewer* atoms through? This is the old, deep mystery of quantum mechanics—the interference of amplitudes. It's the same kind of thing we first saw in the two-slit interference experiment with electrons. We saw that we could get fewer electrons at some places with both slits open than we got with one slit open. It works quantitatively this way. We can write the amplitude that an atom will get through T and S' in the apparatus of (5.17) as the sum of three amplitudes, one for each of the three beams in T ; the sum is equal to zero:

$$\langle 0 S | +T \rangle \langle +T | +S \rangle + \langle 0 S | 0T \rangle \langle 0T | +S \rangle + \langle 0 S | -T \rangle \langle -T | +S \rangle = 0. \quad (5.18)$$

None of the three individual amplitudes is zero—for example, the absolute square of the second amplitude is γ_α , see (5.15)—but the *sum is zero*. We would have also the same answer if S' were set to select the $(-S)$ state. However, in the setup of (5.16), the answer is different. If we call a the amplitude to get through T and S' , in this case we have†

$$a = \langle +S | +T \rangle \langle +T | +S \rangle + \langle +S | 0T \rangle \langle 0T | +S \rangle + \langle +S | -T \rangle \langle -T | +S \rangle = 1. \quad (5.19)$$

In the experiment (5.16) the beam has been split and recombined. Humpty Dumpty has been put back together again. The information about the original $(+S)$ state is retained—it is just as though the T apparatus were not there at all. This is true whatever is put after the “wide-open” T apparatus. We could follow it with an R filter—a filter at some odd angle—or anything we want. The answer will always be the same as if the atoms were taken directly from the first S filter.

So this is the important principle: A T filter—or any filter—with wide-open masks produces no change at all. We should make one additional condition. The wide-open filter must not only transmit all three beams, but it must also *not* produce unequal disturbances on the three beams. For instance, it should not have a strong electric field near one beam and not the others. The reason is that even if this extra disturbance would still let all the atoms through the filter, it could change the *phases* of some of the amplitudes. Then the interference would be changed, and the amplitudes in Eqs. (5.18) and (5.19) would be different. We will always assume that there are no such extra disturbances.

Let's rewrite Eqs. (5.18) and (5.19) in an improved notation. We will let i stand for any one of the three states $(+T)$, $(0T)$, or $(-T)$; then the equations can be written:

$$\sum_{\text{all } i} \langle 0 S | i \rangle \langle i | +S \rangle = 0 \quad (5.20)$$

and

$$\sum_{\text{all } i} \langle +S | i \rangle \langle i | +S \rangle = 1. \quad (5.21)$$

Similarly, for an experiment where S' is replaced by a completely arbitrary filter R , we have

$$\left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_S \quad \left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_T \quad \left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_R. \quad (5.22)$$

† We really cannot conclude from the experiment that $a = 1$, but only that $|a|^2 = 1$, so a might be $e^{i\delta}$, but it can be shown that the choice $\delta = 0$ represents no real loss of generality.

The results will always be the same as if the T apparatus were left out and we had only

$$\left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_S \quad \left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_R.$$

Or, expressed mathematically,

$$\sum_{\text{all } i} \langle +R | i \rangle \langle i | +S \rangle = \langle +R | +S \rangle. \quad (5.23)$$

This is our fundamental law, and it is generally true so long as i stands for the three base states of any filter.

You will notice that in the experiment (5.22) there is no special relation of S and R to T . Furthermore, the arguments would be the same no matter what states they selected. To write the equation in a general way, without having to refer to the specific states selected by S and R , let's call ϕ ("phi") the state prepared by the first filter (in our special example, $+S$) and χ ("khi") the state tested by the final filter (in our example, $+R$). Then we can state our fundamental law of Eq. (5.23) in the form

$$\langle \chi | \phi \rangle = \sum_{\text{all } i} \langle \chi | i \rangle \langle i | \phi \rangle, \quad (5.24)$$

where i is to range over the three base states of some particular filter.

We want to emphasize again what we mean by base states. They are like the three states which can be selected by one of our Stern-Gerlach apparatuses. One condition is that if you have a base state, then the future is independent of the past. Another condition is that if you have a complete set of base states, Eq. (5.24) is true for any set of beginning and ending states ϕ and χ . There is, however, *no unique* set of base states. We began by considering base states *with respect to* a particular apparatus T . We could equally well consider a *different set* of base states with respect to an apparatus S , or with respect to R , etc.† We usually speak of the base states "in a certain representation."

Another condition on a set of base states in any particular representation is that they are all completely different. By that we mean that if we have a $(+T)$ state, there is no amplitude for it to go into a $(0T)$ or a $(-T)$ state. If we let i and j stand for any two base states of a particular set, the general rules discussed in connection with (5.8) are that

$$\langle j | i \rangle = 0$$

for all i and j that are not equal. Of course, we know that

$$\langle i | i \rangle = 1.$$

These two equations are usually written as

$$\langle j | i \rangle = \delta_{ji}, \quad (5.25)$$

where δ_{ji} (the "Kronecker delta") is a symbol that is defined to be zero for $i \neq j$, and to be one for $i = j$.

Equation (5.25) is not independent of the other laws we have mentioned. It happens that we are not particularly interested in the mathematical problem of finding the minimum set of independent axioms that will give all the laws as consequences.‡ We are satisfied if we have a set that is complete and not apparently inconsistent. We can, however, show that Eqs. (5.25) and (5.24) are not independent. Suppose we let ϕ in Eq. (5.24) represent one of the base states of the

† In fact, for atomic systems with three or more base states, there exist other kinds of filters—quite different from a Stern-Gerlach apparatus—which can be used to get more choices for the set of base states (each set with the same *number* of states).

‡ Redundant *truth* doesn't bother us!

same set as i , say the j th state; then we have

$$\langle \chi | j \rangle = \sum_i \langle \chi | i \rangle \langle i | j \rangle.$$

But Eq. (5.25) says that $\langle i | j \rangle$ is zero unless $i = j$, so the sum becomes just $\langle \chi | j \rangle$ and we have an identity, which shows that the two laws are not independent.

We can see that there must be another relation among the amplitudes if both Eqs. (5.10) and (5.24) are true. Equation (5.10) is

$$\langle +T | +S \rangle \langle +T | +S \rangle^* + \langle 0T | +S \rangle \langle 0T | +S \rangle^* + \langle -T | +S \rangle \langle -T | +S \rangle^* = 1.$$

If we write Eq. (5.24), letting both ϕ and χ be the state $(+S)$, the left-hand side is $\langle +S | +S \rangle$, which is clearly $= 1$; so we get once more Eq. (5.19),

$$\langle +S | +T \rangle \langle +T | +S \rangle + \langle +S | 0T \rangle \langle 0T | +S \rangle + \langle +S | -T \rangle \langle -T | +S \rangle = 1.$$

These two equations are consistent (for all relative orientations of the T and S apparatuses) only if

$$\begin{aligned} \langle +S | +T \rangle &= \langle +T | +S \rangle^*, \\ \langle +S | 0T \rangle &= \langle 0T | +S \rangle^*, \\ \langle +S | -T \rangle &= \langle -T | +S \rangle^*. \end{aligned}$$

And it follows that for any states ϕ and χ ,

$$\langle \phi | \chi \rangle = \langle \chi | \phi \rangle^*. \quad (5.26)$$

If this were not true, probability wouldn't be "conserved," and particles would get "lost."

Before going on, we want to summarize the three important general laws about amplitudes. They are Eqs. (5.24), (5.25), and (5.26):

$$\begin{aligned} \text{I} \quad \langle j | i \rangle &= \delta_{ji}, \\ \text{II} \quad \langle \chi | \phi \rangle &= \sum_{\text{all } i} \langle \chi | i \rangle \langle i | \phi \rangle, \\ \text{III} \quad \langle \phi | \chi \rangle &= \langle \chi | \phi \rangle^*. \end{aligned} \quad (5.27)$$

In these equations the i and j refer to *all* the base states of some *one* representation, while ϕ and χ represent any possible states of the atom. It is important to note that II is valid only if the sum is carried out over *all* the base states of the system (in our case, three: $+T$, $0T$, $-T$). These laws say nothing about what we should choose for a base for our set of base states. We began by using a T apparatus, which is a Stern-Gerlach experiment with some arbitrary orientation; but any other orientation, say W , would be just as good. We would have a different set of states to use for i and j , but all the laws would still be good—there is no unique set. One of the great games of quantum mechanics is to make use of the fact that things can be calculated in more than one way.

5-6 The machinery of quantum mechanics

We want to show you why these laws are useful. Suppose we have an atom in a given condition (by which we mean that it was prepared in a certain way), and we want to know what will happen to it in some experiment. In other words, we start with our atom in the state ϕ and want to know what are the *odds* that it will go through some apparatus which accepts atoms only in the condition χ . The laws say that we can describe the apparatus completely in terms of three complex numbers $\langle \chi | i \rangle$, the amplitudes for each base state to be in the condition χ ; and that we can tell what will happen if an atom is put into the apparatus if we describe the state of the atom by giving three numbers $\langle i | \phi \rangle$, the amplitudes for the atom in its original condition to be found in each of the three base states. This is an important idea.

Let's consider another illustration. Think of the following problem: We start with an S apparatus; then we have a complicated mess of junk, which we can call A , and then an R apparatus—like this:

$$\left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\} \left\{ \begin{array}{c} \\ A \\ \end{array} \right\} \left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}. \quad (5.28)$$

$S \qquad \qquad R$

By A we mean any complicated arrangement of Stern-Gerlach apparatuses with masks or half-masks, oriented at peculiar angles, with odd electric and magnetic fields . . . almost anything you want to put. (It's nice to do thought experiments—you don't have to go to all the trouble of actually *building* the apparatus!) The problem then is: With what amplitude does a particle that enters the section A in a $(+S)$ state come out of A in the $(0R)$ state, so that it will get through the last R filter? There is a regular notation for such an amplitude; it is

$$\langle 0R | A | +S \rangle.$$

As usual, it is to be read from right to left (like Hebrew):

$$\langle \text{finish} | \text{through} | \text{start} \rangle.$$

If by chance A doesn't do anything—but is just an open channel—then we write

$$\langle 0R | 1 | +S \rangle = \langle 0R | +S \rangle; \quad (5.29)$$

the two symbols are equivalent. For a more general problem, we might replace $(+S)$ by a general starting state ϕ and $(0R)$ by a general finishing state χ , and we would want to know the amplitude

$$\langle \chi | A | \phi \rangle.$$

A complete analysis of the apparatus A would have to give the amplitude $\langle \chi | A | \phi \rangle$ for every possible pair of states ϕ and χ —an infinite number of combinations! How then can we give a concise description of the behavior of the apparatus A ? We can do it in the following way. Imagine that the apparatus of (5.28) is modified to be

$$\left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\} \left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\} \left\{ \begin{array}{c} \\ A \\ \end{array} \right\} \left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\} \left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}. \quad (5.30)$$

$S \qquad T \qquad \qquad T \qquad R$

This is really no modification at all since the wide-open T apparatuses don't do anything. But they do suggest how we can analyze the problem. There is a certain set of amplitudes $\langle i | +S \rangle$ that the atoms from S will get into the i state of T . Then there is another set of amplitudes that an i state (with respect to T) entering A will come out as a j state (with respect to T). And finally there is an amplitude that each j state will get through the last filter as a $(0R)$ state. For each possible alternative path, there is an amplitude of the form

$$\langle 0R | j \rangle \langle j | A | i \rangle \langle i | +S \rangle,$$

and the total amplitude is the sum of the terms we can get with all possible combinations of i and j . The amplitude we want is

$$\sum_{ij} \langle 0R | j \rangle \langle j | A | i \rangle \langle i | +S \rangle. \quad (5.31)$$

If $(0R)$ and $(+S)$ are replaced by general states χ and ϕ , we would have the same kind of expression; so we have the general result

$$\langle \chi | A | \phi \rangle = \sum_{ij} \langle \chi | j \rangle \langle j | A | i \rangle \langle i | \phi \rangle. \quad (5.32)$$

This then is the machinery of quantum mechanics for a spin-one particle. Every *state* is described by three numbers which are the amplitudes to be in each of some selected set of base states. Every apparatus is described by nine numbers which are the amplitudes to go from one base state to another in the apparatus. From these numbers anything can be calculated.

$$\begin{array}{ccccc}
& & & \text{from} & \\
& & + & 0 & - \\
\text{to} & + & \langle + | A | + \rangle & \langle + | A | 0 \rangle & \langle + | A | - \rangle \\
& 0 & \langle 0 | A | + \rangle & \langle 0 | A | 0 \rangle & \langle 0 | A | - \rangle \\
& - & \langle - | A | + \rangle & \langle - | A | 0 \rangle & \langle - | A | - \rangle
\end{array} \quad (5.33)$$
$$\left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_S \quad \left\{ C \right\} \quad \left\{ \begin{array}{c} + \\ 0 \\ - \end{array} \right\}_R. \quad (5.34)$$
$$\left\{ C \right\} = \left\{ A \right\} \cdot \left\{ B \right\} . \quad (5.35)$$
$$\langle \chi | C | \phi \rangle$$
$$\langle \chi | C | \phi \rangle = \sum_k \langle \chi | B | k \rangle \langle k | A | \phi \rangle.$$
$$\langle j | C | i \rangle = \sum_k \langle j | B | k \rangle \langle k | A | i \rangle. \quad (5.36)$$

5-14

5-7 Transforming to a different base

We want to make one final point about the base states used in the calculations. Suppose we have chosen to work with some particular base—say the S base—and another fellow decides to do the same calculations with a different base—say the T base. To keep things straight let's call our base states the (iS) states, where $i = +, 0, -$. Similarly, we can call his base states (jT) . How can we compare our work with his? The final answers for the result of any measurement should come out the same, but in the calculations the various amplitudes and matrices used will be different. How are they related? For instance, if we both start with the same ϕ , we will describe it in terms of the three amplitudes $\langle iS | \phi \rangle$ that ϕ goes into our base states in the S representation, whereas he will describe it by the amplitudes $\langle jT | \phi \rangle$ that the state ϕ goes into the base states in his T representation. How can we check that we are really both describing the same state ϕ ? We can do it with the general rule II in (5.27). Replacing x by any one of his states jT , we have

$$\langle jT | \phi \rangle = \sum_i \langle jT | iS \rangle \langle iS | \phi \rangle. \quad (5.37)$$

To relate the two representations, we need only give the nine complex numbers of the matrix $\langle jT | iS \rangle$. This matrix can then be used to convert all of his equations to our form. It tells us how to *transform* from one set of base states to another. (For this reason $\langle jT | iS \rangle$ is sometimes called “the transformation matrix from representation S to representation T .” Big words!) Big words!

For the case of spin-one particles for which we have only three base states (for higher spins, there are more) the mathematical situation is analogous to what we have seen in vector algebra. Every vector can be represented by giving three numbers—the components along the axes x , y , and z . That is, every vector can be resolved into three “base” vectors which are vectors along the three axes. But suppose someone else chooses to use a different set of axes— x' , y' , and z' . He will be using different numbers to represent any particular vector. His calculations will look different, but the final results will be the same. We have considered this before and know the rules for transforming vectors from one set of axes to another.

You may want to see how the quantum mechanical transformations work by trying some out; so we will give here, without proof, the transformation matrices for converting the spin-one amplitudes in one representation S to another representation T , for various special relative orientations of the S and T filters. (We will show you in a later chapter how to derive these same results.)

First case: The T apparatus has the same y -axis (along which the particles move) as the S apparatus, but is rotated about the common y -axis by the angle α (as in Fig. 5-6). (To be specific, a set of coordinates x' , y' , z' is fixed in the T apparatus, related to the x , y , z coordinates of the S apparatus by: $z' = z \cos \alpha + x \sin \alpha$, $x' = x \cos \alpha - z \sin \alpha$, $y' = y$.) Then the transformation amplitudes are:

$$\begin{aligned} \langle +T | +S \rangle &= \frac{1}{2}(1 + \cos \alpha), \\ \langle 0T | +S \rangle &= -\frac{1}{\sqrt{2}} \sin \alpha, \\ \langle -T | +S \rangle &= \frac{1}{2}(1 - \cos \alpha), \\ \langle +T | 0S \rangle &= +\frac{1}{\sqrt{2}} \sin \alpha, \\ \langle 0T | 0S \rangle &= \cos \alpha, \\ \langle -T | 0S \rangle &= -\frac{1}{\sqrt{2}} \sin \alpha, \\ \langle +T | -S \rangle &= \frac{1}{2}(1 - \cos \alpha), \\ \langle 0T | -S \rangle &= +\frac{1}{\sqrt{2}} \sin \alpha, \\ \langle -T | -S \rangle &= \frac{1}{2}(1 + \cos \alpha). \end{aligned} \quad (5.38)$$

Second Case: The T apparatus has the same z -axis as S , but is rotated around the z -axis by the angle β . (The coordinate transformation is $z' = z$, $x' = x \cos \beta + y \sin \beta$, $y' = y \cos \beta - x \sin \beta$.) Then the transformation amplitudes are:

$$\begin{aligned}\langle +T | +S \rangle &= e^{+i\beta}, \\ \langle 0T | 0S \rangle &= 1, \\ \langle -T | -S \rangle &= e^{-i\beta}, \\ \text{all others} &= 0.\end{aligned}\tag{5.39}$$

Note that any rotations of T whatever can be made up of the two rotations described.

If a state ϕ is defined by the three numbers

$$C_+ = \langle +S | \phi \rangle, \quad C_0 = \langle 0S | \phi \rangle, \quad C_- = \langle -S | \phi \rangle, \tag{5.40}$$

and the same state is described from the point of view of T by the three numbers

$$C'_+ = \langle +T | \phi \rangle, \quad C'_0 = \langle 0T | \phi \rangle, \quad C'_- = \langle -T | \phi \rangle, \tag{5.41}$$

then the coefficients $\langle jT | iS \rangle$ of (5.38) or (5.39) give the transformation connecting C'_i and C_i . In other words, the C_i are very much like the components of a vector that appear different from the point of view of S and T .

For a spin-one particle *only*—because it requires *three* amplitudes—the correspondence with a vector is very close. In each case, there are three numbers that must transform with coordinate changes in a certain definite way. In fact, there is a set of base states *which transform just like the three components of a vector*. The three combinations

$$C_x = -\frac{1}{\sqrt{2}}(C_+ - C_-), \quad C_y = -\frac{i}{\sqrt{2}}(C_+ + C_-), \quad C_z = C_0 \tag{5.42}$$

transform to C'_x , C'_y , and C'_z just the way that x , y , z transform to x' , y' , z' . [You can check that this is so by using the transformation laws (5.38) and (5.39).] Now you see why a spin-one particle is often called a “vector particle.”

5-8 Other situations

We began by pointing out that our discussion of spin-one particles would be a prototype for any quantum mechanical problem. The generalization has only to do with the numbers of states. Instead of only three base states, any particular situation may involve n base states.† Our basic laws in Eq. (5.27) have exactly the same form—with the understanding that i and j must range over all n base states. Any phenomenon can be analyzed by giving the amplitudes that it starts in each one of the base states and ends in any other one of the base states, and then summing over the complete set of base states. Any proper set of base states can be used, and if someone wishes to use a different set, it is just as good; the two can be connected by using an n by n transformation matrix. We will have more to say later about such transformations.

Finally, we promised to remark on what to do if atoms come directly from a furnace, go through some apparatus, say A , and are then analyzed by a filter which selects the state x . You do not know what the state ϕ is that they start out in. It is perhaps best if you don't worry about this problem just yet, but instead concentrate on problems that always start out with pure states. But if you insist, here is how the problem can be handled.

First, you have to be able to make some reasonable guess about the way the states are distributed in the atoms that come from the furnace. For example, if

† The number of base states n may be, and generally is, infinite.

there were nothing “special” about the furnace, you might reasonably guess that atoms would leave the furnace with random “orientations.” Quantum mechanically, that corresponds to saying that you don’t know anything about the states, but that one-third are in the $(+S)$ state, one-third are in the $(0\ S)$ state, and one-third are in the $(-S)$ state. For those that are in the $(+S)$ state the amplitude to get through is $\langle \chi | A | +S \rangle$ and the probability is $|\langle \chi | A | +S \rangle|^2$, and similarly for the others. The overall probability is then

$$\frac{1}{3}|\langle \chi | A | +S \rangle|^2 + \frac{1}{3}|\langle \chi | A | 0\ S \rangle|^2 + \frac{1}{3}|\langle \chi | A | -S \rangle|^2.$$

Why did we use S rather than, say, T ? The answer is, surprisingly, the same no matter what we choose for our initial resolution—so long as we are dealing with completely random orientations. It comes about in the same way that

$$\sum_i |\langle \chi | iS \rangle|^2 = \sum_j |\langle \chi | jT \rangle|^2$$

for any χ . (We leave it for you to prove.)

Note that it is *not* correct to say that the input state has the amplitudes $\sqrt{1/3}$ to be in $(+S)$, $\sqrt{1/3}$ to be in $(0\ S)$, and $\sqrt{1/3}$ to be in $(-S)$; that would imply that certain interferences might be possible. It is simply that you do not *know* what the initial state is; you have to think in terms of the probability that the system starts out in the various possible initial states, and then you have to take a weighted average over the various possibilities.

Spin One-Half†

6-1 Transforming amplitudes

In the last chapter, using a system of spin one as an example, we outlined the general principles of quantum mechanics:

Any state ψ can be described in terms of a set of base states by giving the amplitudes to be in each of the base states.

The amplitude to go from any state to another can, in general, be written as a sum of products, each product being the amplitude to go into one of the base states times the amplitude to go from that base state to the final condition, with the sum including a term for each base state:

$$\langle \chi | \psi \rangle = \sum_i \langle \chi | i \rangle \langle i | \psi \rangle. \quad (6.1)$$

The base states are orthogonal—the amplitude to be in one if you are in the other is zero:

$$\langle i | j \rangle = \delta_{ij}. \quad (6.2)$$

The amplitude to get from one state to another directly is the complex conjugate of the reverse:

$$\langle \chi | \psi \rangle^* = \langle \psi | \chi \rangle. \quad (6.3)$$

We also discussed a little bit about the fact that there can be more than one base for the states and that we can use Eq. (6.1) to convert from one base to another. Suppose, for example, that we have the amplitudes $\langle iS | \psi \rangle$ to find the state ψ in every one of the base states i of a base system S , but that we then decide that we would prefer to describe the state in terms of another set of base states, say the states j belonging to the base T . In the general formula, Eq. (6.1), we could substitute jT for χ and obtain this formula:

$$\langle jT | \psi \rangle = \sum_i \langle jT | iS \rangle \langle iS | \psi \rangle. \quad (6.4)$$

The amplitudes for the state (ψ) to be in the base states (iT) are related to the amplitudes to be in the base states (iS) by the set of coefficients $\langle jT | iS \rangle$. If there are N base states, there are N^2 such coefficients. Such a set of coefficients is often called the “*transformation matrix*” to go from the *S-representation* to the *T-representation*. This looks rather formidable mathematically, but with a little renaming we can see that it is really not so bad. If we call C_i the amplitude that the state ψ is in the base state iS —that is, $C_i = \langle iS | \psi \rangle$ —and call C'_j the corresponding amplitudes for the base system T —that is, $C'_j = \langle jT | \psi \rangle$, then Eq. (6.4) can be written as

$$C'_j = \sum_i R_{ji} C_i, \quad (6.5)$$

where R_{ji} means the same thing as $\langle jT | iS \rangle$. Each amplitude C'_j is equal to a sum

† This chapter is a rather long and abstract side tour, and it does not introduce any idea which we will not also come to by a different route in later chapters. You can, therefore, skip over it, and come back later if you are interested.

6-1 Transforming amplitudes

6-2 Transforming to a rotated coordinate system

6-3 Rotations about the z-axis

6-4 Rotations of 180° and 90° about y

6-5 Rotations about x

6-6 Arbitrary rotations

over all i of one of the coefficients R_{ji} times each amplitude C_i . It has the same form as the transformation of a vector from one coordinate system to another.

In order to avoid being too abstract for too long, we have given you some examples of these coefficients for the spin-one case, so you can see how to use them in practice. On the other hand, there is a very beautiful thing in quantum mechanics—that from the sheer fact that there are three states and from the symmetry properties of space under rotations, these coefficients can be found purely by abstract reasoning. Showing you such arguments at this early stage has a disadvantage in that you are immersed in another set of abstractions before we get “down to earth.” However, the thing is so beautiful that we are going to do it anyway.

We will show you in this chapter how the transformation coefficients can be derived for spin one-half particles. We pick this case, rather than spin one, because it is somewhat easier. Our problem is to determine the coefficients R_{ji} for a particle—an atomic system—which is split into two beams in a Stern-Gerlach apparatus. We are going to derive all the coefficients for the transformation from one representation to another by pure reasoning—plus a few assumptions. *Some* assumptions are always necessary in order to use “pure” reasoning! Although the arguments will be abstract and somewhat involved, the result we get will be relatively simple to state and easy to understand—and the result is the most important thing. You may, if you wish, consider this as a sort of cultural excursion. We have, in fact, arranged that all the essential results derived here are also derived in some other way when they are needed in later chapters. So you need have no fear of losing the thread of our study of quantum mechanics if you omit this chapter entirely, or study it at some later time. The excursion is “cultural” in the sense that it is intended to show that the principles of quantum mechanics are not only interesting, but are so deep that by adding only a few extra hypotheses about the structure of space, we can deduce a great many properties of physical systems. Also, it is important that we know where the different consequences of quantum mechanics come from, because so long as our laws of physics are incomplete—as we know they are—it is interesting to find out whether the places where our theories fail to agree with experiment is where our logic is the best or where our logic is the worst. Until now, it appears that where our logic is the most abstract it always gives correct results—it agrees with experiment. Only when we try to make specific models of the internal machinery of the fundamental particles and their interactions are we unable to find a theory that agrees with experiment. The theory then that we are about to describe agrees with experiment wherever it has been tested—for the strange particles as well as for electrons, protons, and so on.

One remark on an annoying, but interesting, point before we proceed: It is not possible to determine the coefficients R_{ji} uniquely, because there is always some arbitrariness in the probability amplitudes. If you have a set of amplitudes of any kind, say the amplitudes to arrive at some place by a whole lot of different routes, and if you multiply every single amplitude by the same phase factor—say by $e^{i\delta}$ —you have another set that is just as good. So, it is always possible to make an arbitrary change in phase of all the amplitudes in any given problem if you want to.

Suppose you calculate some probability by writing a sum of several amplitudes, say $(A + B + C + \cdots)$ and taking the absolute square. Then somebody else calculates the same thing by using the sum of the amplitudes $(A' + B' + C' + \cdots)$ and taking the absolute square. If all the A' , B' , C' , etc., are equal to the A , B , C , etc., except for a factor $e^{i\delta}$, all probabilities obtained by taking the absolute squares will be exactly the same, since $(A' + B' + C' + \cdots)$ is then equal to $e^{i\delta}(A + B + C + \cdots)$. Or suppose, for instance, that we were computing something with Eq. (6.1), but then we suddenly change all of the phases of a certain base system. Every one of the amplitudes $\langle i | \psi \rangle$ would be multiplied by the same factor $e^{i\delta}$. Similarly, the amplitudes $\langle i | \chi \rangle$ would also be changed by $e^{i\delta}$, but the amplitudes $\langle \chi | i \rangle$ are the complex conjugates of the amplitudes $\langle i | \chi \rangle$; therefore, the former gets changed by the factor $e^{-i\delta}$. The plus and minus $i\delta$'s

in the exponents cancel out, and we would have the same expression we had before. So it is a general rule that if we change all the amplitudes with respect to a given base system by the same phase—or even if we just change *all* the amplitudes in any problem by the same phase—it makes no difference. There is, therefore, some freedom to choose the phases in our transformation matrix. Every now and then we will make such an arbitrary choice—usually following the conventions that are in general use.

6-2 Transforming to a rotated coordinate system

We consider again the “improved” Stern-Gerlach apparatus described in the last chapter. A beam of spin one-half particles, entering at the left, would, in general, be split into *two* beams, as shown schematically in Fig. 6-1. (There were three beams for spin *one*.) As before, the beams are put back together again unless one or the other of them is blocked off by a “stop” which intercepts the beam at its half-way point. In the figure we show an arrow which points in the direction of the increase of the *magnitude* of the field—say toward the magnet pole with the sharp edges. This arrow we take to represent the “up” axis of any particular apparatus. It is fixed relative to the apparatus and will allow us to indicate the relative orientations when we use several apparatuses together. We also assume that the direction of the magnetic field in each magnet is always the same with respect to the arrow.

We will say that those atoms which go in the “upper” beam are in the (+) state *with respect to that apparatus* and that those in the “lower” beam are in the (−) state. (There is no “zero” state for spin one-half particles.)

Now suppose we put two of our modified Stern-Gerlach apparatuses in sequence, as shown in Fig. 6-2(a). The first one, which we call *S*, can be used to prepare a pure (+*S*) or a pure (−*S*) state by blocking one beam or the other. [As shown it prepares a pure (+*S*) state.] For each condition, there is some amplitude for a particle that comes out of *S* to be in either the (+*T*) or the (−*T*) beam of the second apparatus. There are, in fact, just four amplitudes: the amplitude to go from (+*S*) to (+*T*), from (+*S*) to (−*T*), from (−*S*) to (+*T*), from (−*S*) to (−*T*). These amplitudes are just the four coefficients of the transformation matrix R_{ji} to go from the *S*-representation to the *T*-representation. We can consider that the first apparatus “prepares” a particular state in one representation and that the second apparatus “analyzes” that state in terms of the second representation. The kind of question we want to answer, then, is this: If an atom has been prepared in a given condition—say the (+*S*) state—by blocking one of the beams in the apparatus *S*, what is the chance that it will get through the second apparatus *T* if this is set for, say, the (−*T*) state. The result will depend, of course, on the angles between the two systems *S* and *T*.

We should explain why it is that we could have any hope of finding the coefficients R_{ji} by deduction. You know that it is almost impossible to believe that if a particle has its spin lined up in the +*z*-direction, that there is some chance of finding the same particle with its spin pointing in the +*x*-direction—or in any other direction at all. In fact, it *is* almost impossible, but not quite. It is so nearly impossible that there is *only one way* it can be done, and that is the reason we can find out what that unique way is.

The first kind of argument we can make is this. Suppose we have a setup like the one in Fig. 6-2(a), in which we have the two apparatuses *S* and *T*, with *T* cocked at the angle α with respect to *S*, and we let only the (+) beam through *S* and the (−) beam through *T*. We would observe a certain number for the probability that the particles coming out of *S* get through *T*. Now suppose we make another measurement with the apparatus of Fig. 6-2(b). The *relative* orientation of *S* and *T* is the same, but the whole system sits at a different angle in space. We want to *assume* that both of these experiments give the same number for the chance that a particle in a pure state with respect to *S* will get into some particular state with respect to *T*. We are assuming, in other words, that the result of any experiment of this type is the same—that the *physics* is the same—no matter

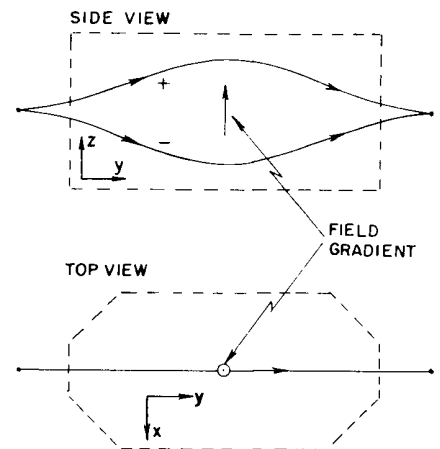


Fig. 6-1. Top and side views of an “improved” Stern-Gerlach apparatus with beams of a spin one-half particle.

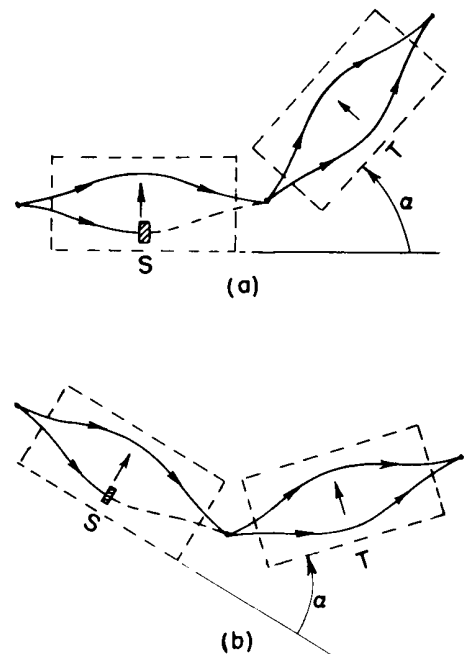


Fig. 6-2. Two equivalent experiments.

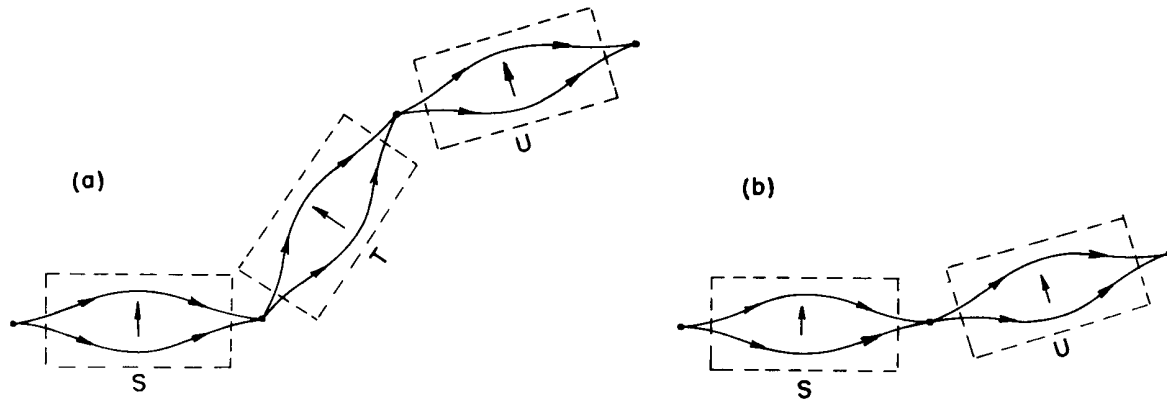


Fig. 6-3. If T is "wide open," (b) is equivalent to (a).

how the *whole* apparatus is oriented in space. (You say, "That's obvious." But it *is* an assumption, and it is "right" only if it is actually what happens.) That means that the coefficients R_{ji} depend only on the relation in space of S and T , and not on the absolute situation of S and T . To say this in another way, R_{ji} depends only on the *rotation* which carries S to T , for evidently what is the same in Fig. 6-2(a) and Fig. 6-2(b) is the three-dimensional rotation which would carry apparatus S into the orientation of apparatus T . When the transformation matrix R_{ji} depends only on a rotation, as it does here, it is called a *rotation matrix*.

For our next step we will need one more piece of information. Suppose we add a third apparatus which we can call U , which follows T at some arbitrary angle, as in Fig. 6-3(a). (It's beginning to look horrible, but that's the fun of abstract thinking—you can make the most weird experiments just by drawing lines!) Now what is the $S \rightarrow T \rightarrow U$ transformation? What we really want to ask for is the amplitude to go from some state with respect to S to some other state with respect to U , when we know the transformation from S to T and from T to U . We are then asking about an experiment in which both channels of T are open. We can get the answer by applying Eq. (6.5) twice in succession. For going from the S -representation to the T -representation, we have

$$C'_j = \sum_i R_{ji}^{TS} C_i, \quad (6.6)$$

where we put the superscripts TS on the R , so that we can distinguish it from the coefficients R^{UT} we will have for going from T to U .

Assuming the amplitudes to be in the base states of the U -representation C''_k , we can relate them to the T -amplitudes by using Eq. (6.5) once more; we get

$$C''_k = \sum_j R_{kj}^{UT} C'_j. \quad (6.7)$$

Now we can combine Eqs. (6.6) and (6.7) to get the transformation to U directly from S . Substituting C'_j from Eq. (6.6) in Eq. (6.7), we have

$$C''_k = \sum_j R_{kj}^{UT} \sum_i R_{ji}^{TS} C_i. \quad (6.8)$$

Or, since i does not appear in R_{kj}^{UT} , we can put the i -summation also in front, and write

$$C''_k = \sum_i \sum_j R_{kj}^{UT} R_{ji}^{TS} C_i. \quad (6.9)$$

This is the formula for a double transformation.

Notice, however, that so long as all the beams in T are unblocked, the state coming out of T is the same as the one that went in. We could just as well have made a transformation from the S -representation directly to the U -representation. It should be the same as putting the U apparatus right after S , as in Fig.

6-3(b). In that case, we would have written

$$C_k'' = \sum_i R_{ki}^{US} C_i, \quad (6.10)$$

with the coefficients R_{ki}^{US} belonging to this transformation. Now, clearly, Eqs. (6.9) and (6.10) should give the same amplitudes C_k'' , and this should be true no matter what the original state ϕ was which gave us the amplitudes C_i . So it must be that

$$R_{ki}^{US} = \sum_j R_{kj}^{UT} R_{ji}^{TS}. \quad (6.11)$$

In other words, for any rotation $S \rightarrow U$ of a reference base, which is viewed as a compounding of two successive rotations $S \rightarrow T$ and $T \rightarrow U$, the rotation matrix R_{ki}^{US} can be obtained from the matrices of the two partial rotations by Eq. (6.11). If you wish, you can find Eq. (6.11) directly from Eq. (6.1), for it is only a different notation for $\langle kU | iS \rangle = \sum_j \langle kU | jT \rangle \langle jT | iS \rangle$.

To be thorough, we should add the following parenthetical remarks. They are not terribly important, however, so you can skip to the next section if you want. What we have said is not quite right. We cannot really say that Eq. (6.9) and Eq. (6.10) must give *exactly* the same amplitudes. Only the *physics* should be the same; all the amplitudes could be different by some common phase factor like $e^{i\delta}$ without changing the result of any calculation about the real world. So, instead of Eq. (6.11), all we can say, really, is that

$$e^{i\delta} R_{ki}^{US} = \sum_j R_{kj}^{UT} R_{ji}^{TS}, \quad (6.12)$$

where δ is *some* real constant. What this extra factor of $e^{i\delta}$ means, of course, is that the amplitudes we get if we use the matrix R^{US} might all differ by the same phase ($e^{-i\delta}$) from the amplitude we would get using the two rotations R^{UT} and R^{TS} . We know that it doesn't matter if all amplitudes are changed by the same phase, so we could just ignore this phase factor if we wanted to. It turns out, however, that if we define all of our rotation matrices in a particular way, this extra phase factor will never appear—the δ in Eq. (6.12) will always be zero. Although it is not important for the rest of our arguments, we can give a quick proof by using a mathematical theorem about determinants. [If you don't yet know much about determinants, don't worry about the proof and just skip to the definition of Eq. (6.15).]

First, we should say that Eq. (6.11) is the mathematical definition of a “product” of two matrices. (It is just convenient to be able to say: “ R^{US} is the product of R^{UT} and R^{TS} .”) Second, there is a theorem of mathematics—which you can easily prove for the two-by-two matrices we have here—which says that the determinant of a “product” of two matrices is the product of their determinants. Applying this theorem to Eq. (6.12), we get

$$e^{i2\delta} (\text{Det } R^{US}) = (\text{Det } R^{UT}) \cdot (\text{Det } R^{TS}). \quad (6.13)$$

(We leave off the subscripts, because they don't tell us anything useful.) Yes, the 2δ is right. Remember that we are dealing with two-by-two matrices; every term in the matrix R_{ki}^{US} is multiplied by $e^{i\delta}$, so each product in the determinant—which has *two* factors—gets multiplied by $e^{i2\delta}$. Now let's take the square root of Eq. (6.13) and divide it into Eq. (6.12); we get

$$\frac{R_{ki}^{US}}{\sqrt{\text{Det } R^{US}}} = \sum_j \frac{R_{kj}^{UT}}{\sqrt{\text{Det } R^{UT}}} \frac{R_{ji}^{TS}}{\sqrt{\text{Det } R^{TS}}}. \quad (6.14)$$

The extra phase factor has disappeared.

Now it turns out that if we want all of our amplitudes in any given representation to be normalized (which means, you remember, that $\sum_i \langle \phi | i \rangle \langle i | \phi \rangle = 1$), the rotation matrices will all have determinants that are pure imaginary exponentials, like $e^{i\alpha}$. (We won't prove it; you will see that it always comes out that way.) So we can, if we wish, choose to make all our rotation matrices R have a unique phase by making $\text{Det } R = 1$. It is done like this. Suppose we find a rotation matrix R in some arbitrary way. We make it a rule to “convert” it to “standard form” by defining

$$R_{\text{standard}} = \frac{R}{\sqrt{\text{Det } R}}. \quad (6.15)$$

We can do this because we are just multiplying each term of R by the same phase factor, to get the phases we want. In what follows, we will always assume that our matrices have been put in the “standard form”; then we can use Eq. (6.11) without having any extra phase factors.

6-3 Rotations about the z -axis

We are now ready to find the transformation matrix R_{ji} between two different representations. With our rule for compounding rotations and our assumption that space has no preferred direction, we have the keys we need for finding the matrix of any arbitrary rotation. There is only *one* solution. We begin with the transformation which corresponds to a rotation about the z -axis. Suppose we have two apparatuses S and T placed in series along a straight line with their axes parallel and pointing out of the page, as shown in Fig. 6-4(a). We take our “ z -axis” in this direction. Surely, if the beam goes “up” (toward $+z$) in the S apparatus, it will do the same in the T apparatus. Similarly, if it goes down in S , it will go down in T . Suppose, however, that the T apparatus were placed at some other angle, but still with its axis parallel to the axis of S , as in Fig. 6-4(b). Intuitively, you would say that a $(+)$ beam in S would still go with a $(+)$ beam in T , because the fields and field gradients are still in the same physical direction. And that would be quite right. Also, a $(-)$ beam in S would still go into a $(-)$ beam in T . The same result would apply for any orientation of T in the xy -plane of S . What does this tell us about the relation between $C'_+ = \langle +T | \psi \rangle$, $C'_- = \langle -T | \psi \rangle$ and $C_+ = \langle +S | \psi \rangle$, $C_- = \langle -S | \psi \rangle$? You might conclude that any rotation about the z -axis of the “frame of reference” for base states leaves the amplitudes C_+ to be “up” and “down,” the same as before. We could write $C'_+ = C_+$ and $C'_- = C_-$ —but that is *wrong*. All we *can* conclude is that for such rotations the probabilities to be in the “up” beam are the same for the S and T apparatuses. That is,

$$|C'_+| = |C_+| \quad \text{and} \quad |C'_-| = |C_-|.$$

We cannot say that the *phases* of the amplitudes referred to the T apparatus may not be different for the two different orientations in (a) and (b) of Fig. 6-4.

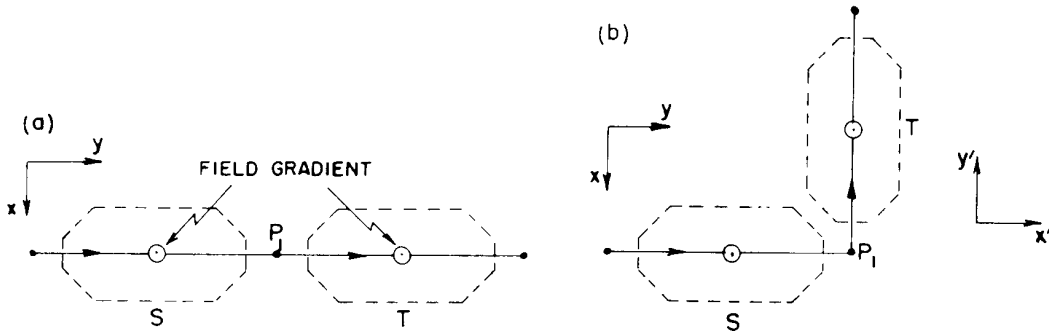


Fig. 6-4. Rotating 90° about the z -axis.

The two apparatuses in (a) and (b) of Fig. 6-4 are, in fact, different, as we can see in the following way. Suppose that we put an apparatus in front of S which produces a pure $(+x)$ state. (The x -axis points toward the bottom of the figure.) Such particles would be split into $(+z)$ and $(-z)$ beams in S , but the two beams would be recombined to give a $(+x)$ state again at P_1 —the exit of S . The same thing happens again in T . If we follow T by a third apparatus U , whose axis is in the $(+x)$ direction and, as shown in Fig. 6-5(a), all the particles would go into the $(+)$ beam of U . Now imagine what happens if T and U are swung around *together* by 90° to the positions shown in Fig. 6-5(b). Again, the T apparatus puts out just what it takes in, so the particles that enter U are in a $(+x)$ state with respect to S . But U now analyzes for the $(+y)$ state with respect to S , which is different. (By symmetry, we would now expect only one-half of the particles to get through.)

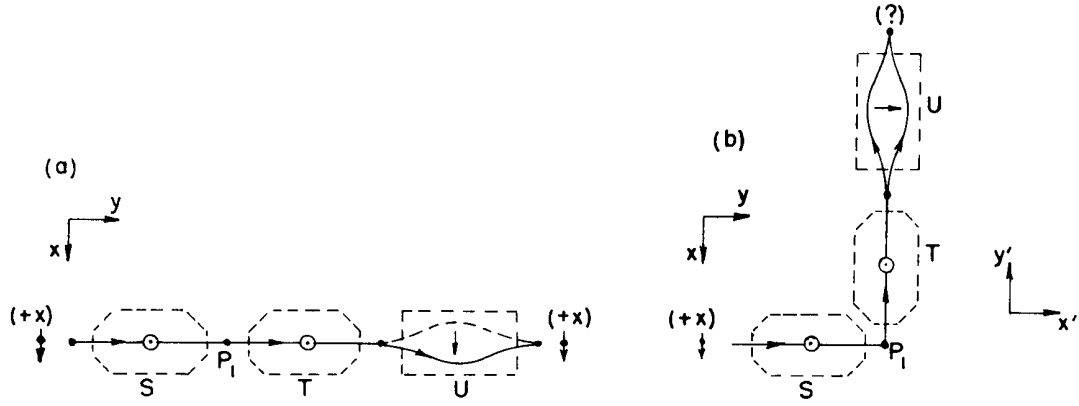


Fig. 6-5. Particle in a $(+x)$ state behaves differently in (a) and (b).

What could have changed? The apparatuses T and U are still in the same *physical* relationship to each other. Can the *physics* be changed just because T and U are in a different orientation? Our original assumption is that it should not. It must be that the *amplitudes* with respect to T are different in the two cases shown in Fig. 6-5—and, therefore, also in Fig. 6-4. There must be some way for a particle to know that it has turned the corner at P_1 . How could it tell? Well, all we have decided is that the *magnitudes* of C'_1 and C'_2 are the same in the two cases, but they could—in fact, *must*—have different *phases*. We conclude that C'_+ and C_+ must be related by

$$C'_+ = e^{i\lambda} C_+,$$

and that C'_- and C_- must be related by

$$C'_- = e^{i\mu} C_-,$$

where λ and μ are real numbers which must be related in some way to the angle between S and T .

The only thing we can say at the moment about λ and μ is that they must not be equal [except for the special case shown in Fig. 6-5(a), when T is in the same orientation as S]. We have seen that equal phase changes in all amplitudes have no physical consequence. For the same reason, we can always add the same arbitrary amount to both λ and μ without changing anything. So we are permitted to *choose* to make λ and μ equal to plus and minus the same number. That is, we can always take

$$\lambda' = \lambda - \frac{(\lambda + \mu)}{2}, \quad \mu' = \mu - \frac{(\lambda + \mu)}{2}.$$

Then

$$\lambda' = \frac{\lambda}{2} - \frac{\mu}{2} = -\mu'.$$

So we adopt the convention† that $\mu = -\lambda$. We have then the general rule that for a rotation of the reference apparatus by some angle about the z -axis, the transformation is

$$C'_+ = e^{+i\lambda} C_+, \quad C'_- = e^{-i\lambda} C_-. \quad (6.16)$$

The absolute values are the same, only the phases are different. These phase factors are responsible for the different results in the two experiments of Fig. 6-5.

Now we would like to know the law that relates λ to the angle between S and T . We already know the answer for one case. If the angle is zero, λ is zero. Now we will *assume* that the phase shift λ is a continuous function of angle ϕ between S and T (see Fig. 6-4) as ϕ goes to zero—as only seems reasonable. In

† Looking at it another way, we are just putting the transformation in the “standard form” described in Section 6-2 by using Eq. (6.15).

ther words, if we rotate T from the straight line through S by the small angle ϵ , the λ is also a small quantity, say $m\epsilon$, where m is some number. We write it this way because we can show that λ must be proportional to ϵ . Suppose we were to put after T another apparatus T' which makes the angle ϵ with T , and, therefore, the angle 2ϵ with S . Then, with respect to T , we have

$$C'_+ = e^{i\lambda} C_+,$$

and with respect to T' , we have

$$C''_+ = e^{i\lambda} C'_+ = e^{i2\lambda} C_+.$$

But we know that we should get the same result if we put T' right after S . Thus, when the angle is doubled, the phase is doubled. We can evidently extend the argument and build up any rotation at all by a sequence of infinitesimal rotations. We conclude that for *any* angle ϕ , λ is proportional to the angle. We can, therefore, write $\lambda = m\phi$.

The general result we get, then, is that for T rotated about the z -axis by the angle ϕ with respect to S

$$C'_+ = e^{im\phi} C_+, \quad C'_- = e^{-im\phi} C_-. \quad (6.17)$$

For the angle ϕ , and for all rotations we speak of in the future, we adopt the standard convention that a *positive* rotation is a *right-handed* rotation about the positive direction of the reference axis. A positive ϕ has the sense of rotation of a right-handed screw advancing in the positive z -direction.

Now we have to find what m must be. First, we might try this argument: Suppose T is rotated by 360° ; then, clearly, it is right back at zero degrees, and we should have $C'_+ = C_+$ and $C'_- = C_-$, or, what is the same thing, $e^{im2\pi} = 1$. We get $m = 1$. *This argument is wrong!* To see that it is, consider that T is rotated by 180° . If m were equal to 1, we would have $C'_+ = e^{i\pi} C_+ = -C_+$ and $C'_- = e^{-i\pi} C_- = -C_-$. However, this is just the *original* state all over again. *Both* amplitudes are just multiplied by -1 which gives back the original physical system. (It is again a case of a common phase change.) This means that if the angle between T and S in Fig. 6-5(b) is increased to 180° , the system (with respect to T) would be indistinguishable from the zero-degree situation, and the particles would again go through the $(+)$ state of the U apparatus. At 180° , though, the $(+)$ state of the U apparatus is the $(-x)$ state of the original S apparatus. So a $(+x)$ state would become a $(-x)$ state. But we have done nothing to *change* the original state; the answer is wrong. We cannot have $m = 1$.

We must have the situation that a rotation by 360° and *no smaller angle* reproduces the same physical state. This will happen if $m = \frac{1}{2}$. Then, and only then, will the first angle that reproduces the same *physical* state be $\phi = 360^\circ$.† It gives

$$\left. \begin{aligned} C'_+ &= -C_+ \\ C'_- &= -C_- \end{aligned} \right\} 360^\circ \text{ about } z\text{-axis}. \quad (6.18)$$

It is very curious to say that if you turn the apparatus 360° you get new amplitudes. They aren't really new, though, because the common change of sign doesn't give any different physics. If someone else had decided to change all the signs of the amplitudes because he thought he had turned 360° , that's all right; he gets the same physics.‡ So our final answer is that if we know the amplitudes C_+ and C_- for spin one-half particles with respect to a reference frame S , and we then use a base

† It appears that $m = -\frac{1}{2}$ would also work. However, we see in (6.17) that the change in sign merely redefines the notation for a spin-up particle.

‡ Also, if something has been rotated by a sequence of small rotations whose net result is to return it to the original orientation, it is possible to define the idea that it has been rotated 360° —as distinct from zero net rotation—if you have kept track of the whole history. (Interestingly enough, this is *not* true for a net rotation of 720° .)

system referred to T which is obtained from S by a rotation of ϕ around the z -axis, the new amplitudes are given in terms of the old by

$$\left. \begin{aligned} C'_+ &= e^{i\phi/2} C_+ \\ C'_- &= e^{-i\phi/2} C_- \end{aligned} \right\} \phi \text{ about } z. \quad (6.19)$$

6-4 Rotations of 180° and 90° about y

Next, we will try to guess the transformation for a rotation of T with respect to S of 180° around an axis *perpendicular* to the z -axis—say, about the y -axis. (We have defined the coordinate axes in Fig. 6-1.) In other words, we start with two identical Stern-Gerlach equipments, with the second one, T , turned “upside down” with respect to the first one, S , as in Fig. 6-6. Now if we think of our particles as little magnetic dipoles, a particle that is the $(+S)$ state—so that it goes on the “upper” path in the first apparatus—will also take the “upper” path in the second, so that it will be in the *minus* state with respect to T . (In the inverted T apparatus, both the gradients *and* the field direction are reversed; for a particle with its magnetic moment in a given direction, the force is unchanged.) Anyway, what is “up” with respect to S will be “down” with respect to T . For these relative positions of S and T , then, we know that the transformation must give

$$|C'_+| = |C_-|, \quad |C'_-| = |C_+|.$$

As before, we cannot rule out some additional phase factors; we could have (for 180° about the y -axis)

$$C'_+ = e^{i\beta} C_- \quad \text{and} \quad C'_- = e^{i\gamma} C_+, \quad (6.20)$$

where β and γ are still to be determined.

What about a rotation of 360° about the y -axis? Well, we already know the answer for a rotation of 360° about the z -axis—the amplitude to be in any state changes sign. A rotation of 360° around any axis always brings us back to the original position. It must be that for *any* 360° rotation, the result is the same as a 360° rotation about the z -axis—all amplitudes simply change sign. Now suppose we imagine two successive rotations of 180° about y —using Eq. (6.20)—we should get the result of Eq. (6.18). In other words,

$$C''_+ = e^{i\beta} C'_- = e^{i\beta} e^{i\gamma} C_+ = -C_+ \quad (6.21)$$

and

$$C''_- = e^{i\gamma} C'_+ = e^{i\gamma} e^{i\beta} C_- = -C_-.$$

This means that

$$e^{i\beta} e^{i\gamma} = -1 \quad \text{or} \quad e^{i\gamma} = -e^{-i\beta}.$$

So the transformation for a rotation of 180° about the y -axis can be written

$$C'_+ = e^{i\beta} C_-, \quad C'_- = -e^{-i\beta} C_+. \quad (6.22)$$

The arguments we have just used would apply equally well to a rotation of 180° about *any* axis in the xy -plane, although different axes can, of course, give different numbers for β . However, that is the only way they can differ. Now there is a certain amount of arbitrariness in the number β , but once it is specified for one axis of rotation in the xy -plane it is determined for any other axis. It is *conventional* to choose to set $\beta = 0$ for a 180° rotation about the y -axis.

To show that we have this choice, suppose we imagine that β was not equal to zero for a rotation about the y -axis; then we can show that there is *some other* axis in the xy -plane, for which the corresponding phase factor *will* be zero. Let's find the phase factor β_A for an axis A that makes the angle α with the y -axis, as shown in Fig. 6-7(a). (For clarity, the figure is drawn with α equal to a negative number, but that doesn't matter.) Now if we take a T apparatus which is initially lined up with the S apparatus and is then rotated 180° about the axis A , its axes—which we will call x'' , y'' , and z'' —will be as shown in Fig. 6-7(a). The amplitudes

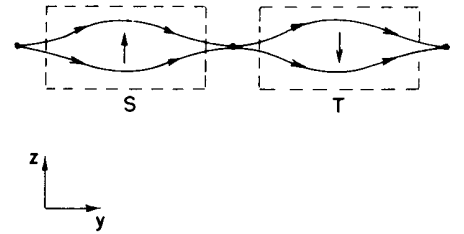


Fig. 6-6. A rotation of 180° about the y -axis.

with respect to T will then be

$$C'_+ = e^{i\beta_A} C_-, \quad C''_- = -e^{-i\beta_A} C_+. \quad (6.23)$$

We can now think of getting to the same orientation by the two successive rotations shown in (b) and (c) of the figure. First, we imagine an apparatus U which is rotated with respect to S by 180° about the y -axis. The axes x' , y' , and z' of U will be as shown in Fig. 6-7(b), and the amplitudes *with respect to U* are given by (6.22).

Now notice that we can go from U to T by a rotation about the “ z -axis” of U , namely about z' , as shown in Fig. 6-7(c). From the figure you can see that the angle required is two times the angle α but in the opposite direction (with respect to z'). Using the transformation of (6.19) with $\phi = -2\alpha$, we get

$$C''_+ = e^{-i\alpha} C'_+, \quad C''_- = e^{+i\alpha} C'_-. \quad (6.24)$$

Combining Eqs. (6.24) and (6.22), we get that

$$C''_+ = e^{i(\beta-\alpha)} C_-, \quad C''_- = -e^{-i(\beta-\alpha)} C_+. \quad (6.25)$$

These amplitudes must, of course, be the same as we got in 6.23). So β_A must be related to α and β by

$$\beta_A = \beta - \alpha. \quad (6.26)$$

This means that if the angle α between the A -axis and the y -axis (of S) is equal to β , the transformation for a rotation of 180° about A will have $\beta_A = 0$.

Now so long as *some* axis perpendicular to the z -axis is going to have $\beta = 0$, we may as well take it to be the y -axis. It is purely a matter of *convention*, and we adopt the one in general use. Our result: For a rotation of 180° about the y -axis, we have

$$\left. \begin{aligned} C'_+ &= C_- \\ C'_- &= -C_+ \end{aligned} \right\} 180^\circ \text{ about } y. \quad (6.27)$$

While we are thinking about the y -axis, let's next ask for the transformation matrix for a rotation of 90° about y . We can find it because we know that two successive 90° rotations about the same axis must equal one 180° rotation. We start by writing the transformation for 90° in the most general form:

$$C'_+ = aC_+ + bC_-, \quad C'_- = cC_+ + dC_-. \quad (6.28)$$

A second rotation of 90° about the same axis would have the same coefficients:

$$C''_+ = aC'_+ + bC'_-, \quad C''_- = cC'_+ + dC'_-. \quad (6.29)$$

Combining Eqs. (6.28) and (6.29), we have

$$\begin{aligned} C''_+ &= a(aC_+ + bC_-) + b(cC_+ + dC_-), \\ C''_- &= c(aC_+ + bC_-) + d(cC_+ + dC_-). \end{aligned} \quad (6.30)$$

However, from (6.27) we know that

$$C''_+ = C_-, \quad C''_- = -C_+,$$

so that we must have that

$$\begin{aligned} ab + bd &= 1, \\ a^2 + bc &= 0, \\ ac + cd &= -1, \\ bc + d^2 &= 0. \end{aligned} \quad (6.31)$$

These four equations are enough to determine all our unknowns: a , b , c , and d .

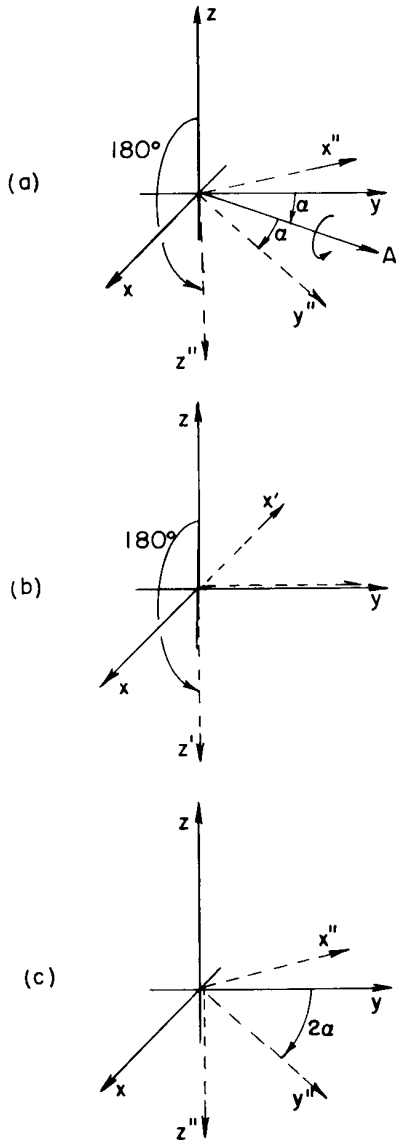


Fig. 6-7. A 180° rotation about the axis A is equivalent to a rotation of 180° about y , followed by a rotation about z' .

It is not hard to do. Look at the second and fourth equations. Deduce that $a^2 = d^2$, which means that $a = d$ or else that $a = -d$. But $a = -d$ is out, because then the first equation wouldn't be right. So $d = a$. Using this, we have immediately that $b = 1/2a$ and that $c = -1/2a$. Now we have everything in terms of a . Putting, say, the second equation all in terms of a , we have

$$a^2 - \frac{1}{4a^2} = 0 \quad \text{or} \quad a^4 = \frac{1}{4}.$$

This equation has four different solutions, but only two of them give the standard value for the determinant. We might as well take $a = 1/\sqrt{2}$; then†

$$\begin{aligned} a &= 1/\sqrt{2}, & b &= 1/\sqrt{2}, \\ c &= -1/\sqrt{2}, & d &= 1/\sqrt{2}. \end{aligned}$$

In other words, for two apparatuses S and T , with T rotated with respect to S by 90° about the y -axis, the transformation is

$$\left. \begin{aligned} C'_+ &= \frac{1}{\sqrt{2}} (C_+ + C_-) \\ C'_- &= \frac{1}{\sqrt{2}} (-C_+ + C_-) \end{aligned} \right\} 90^\circ \text{ about } y. \quad (6.32)$$

We can, of course, solve these equations for C_+ and C_- , which will give us the transformation for a rotation of *minus* 90° about y . Changing the primes around, we would conclude that

$$\left. \begin{aligned} C'_+ &= \frac{1}{\sqrt{2}} (C_+ - C_-) \\ C'_- &= \frac{1}{\sqrt{2}} (C_+ + C_-) \end{aligned} \right\} -90^\circ \text{ about } y. \quad (6.33)$$

6-5 Rotations about x

You may be thinking: "This is getting ridiculous. What are they going to do next, 47° around y , then 33° about x , and so on, forever?" No, we are almost finished. With just two of the transformations we have— 90° about y , and an arbitrary angle about z (which we did first if you remember)—we can generate any rotation at all.

As an illustration, suppose that we want the angle α around x . We know how to deal with the angle α around z , but now we want it around x . How do we get it? First, we turn the axis z down onto x —which is a rotation of $+90^\circ$ about y , as shown in Fig. 6-8. Then we turn through the angle α around z' . Then we rotate -90° about y'' . The net result of the three rotations is the same as turning around x by the angle α . It is a property of space.

(These facts of the combinations of rotations, and what they produce, are hard to grasp intuitively. It is rather strange, because we live in three dimensions, but it is hard for us to appreciate what happens if we turn this way and then that way. Perhaps, if we were fish or birds and had a real appreciation of what happens when we turn somersaults in space, we could more easily appreciate such things.)

Anyway, let's work out the transformation for a rotation by α around the x -axis by using what we know. From the first rotation by $+90^\circ$ around y the amplitudes go according to Eq. (6.32). Calling the rotated axes x' , y' , and z' , the

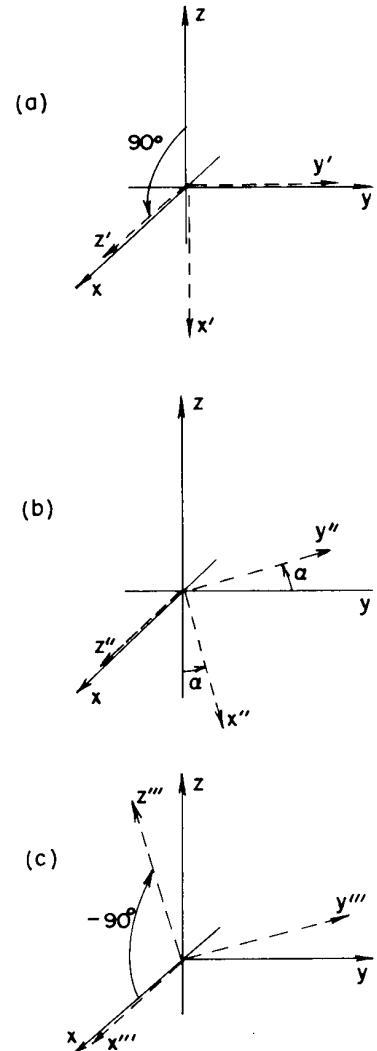


Fig. 6-8. A rotation by α about the x -axis is equivalent to: (a) a rotation by $+90^\circ$ about y , followed by (b) a rotation by α about z' , followed by (c) a rotation of -90° about y'' .

† The other solution changes all signs of a , b , c , and d and corresponds to a -270° rotation.

next rotation by the angle α around z' takes us to a frame x'', y'', z'' , for which

$$C''_+ = e^{i\alpha/2} C'_+, \quad C''_- = e^{-i\alpha/2} C'_-.$$

The last rotation of -90° about y'' takes us to x''', y''', z''' ; by (6.33),

$$C'''_+ = \frac{1}{\sqrt{2}} (C''_+ - C''_-), \quad C'''_- = \frac{1}{\sqrt{2}} (C''_+ + C''_-).$$

Combining these last two transformations, we get

$$C'''_+ = \frac{1}{\sqrt{2}} (e^{+i\alpha/2} C'_+ - e^{-i\alpha/2} C'_-),$$

$$C'''_- = \frac{1}{\sqrt{2}} (e^{+i\alpha/2} C'_+ + e^{-i\alpha/2} C'_-).$$

Using Eqs. (6.32) for C'_+ and C'_- , we get the complete transformation:

$$C'''_+ = \frac{1}{2} \{ e^{+i\alpha/2} (C_+ + C_-) - e^{-i\alpha/2} (-C_+ + C_-) \},$$

$$C'''_- = \frac{1}{2} \{ e^{+i\alpha/2} (C_+ + C_-) + e^{-i\alpha/2} (-C_+ + C_-) \}.$$

We can put these formulas in a simpler form by remembering that

$$e^{i\theta} + e^{-i\theta} = 2 \cos \theta, \quad \text{and} \quad e^{i\theta} - e^{-i\theta} = 2i \sin \theta.$$

We get

$$\left. \begin{aligned} C'''_+ &= \left(\cos \frac{\alpha}{2} \right) C_+ + i \left(\sin \frac{\alpha}{2} \right) C_- \\ C'''_- &= i \left(\sin \frac{\alpha}{2} \right) C_+ + \left(\cos \frac{\alpha}{2} \right) C_- \end{aligned} \right\} \alpha \text{ about } x. \quad (6.34)$$

Here is our transformation for a rotation about the x -axis by *any* angle α . It is only a little more complicated than the others.

6-6 Arbitrary rotations

Now we can see how to do *any* angle at all. First, notice that any relative orientation of two coordinate frames can be described in terms of three angles, as shown in Fig. 6-9. If we have a set of axes $x', y',$ and z' oriented in any way at all with respect to $x, y,$ and z , we can describe the relationship between the two frames by means of the three Euler angles $\alpha, \beta,$ and γ , which define three successive rotations that will bring the x, y, z frame into the x', y', z' frame. Starting at x, y, z , we rotate our frame through the angle β about the z -axis, bringing the x -axis to the line x_1 . Then, we rotate by α about this temporary x -axis, to bring z down to z' . Finally, a rotation about the new z -axis (that is, z') by the angle γ will bring the x -axis into x' and the y -axis into y' .† We know the transformations for each of the three rotations—they are given in (6.19) and (6.34). Combining them in the proper order, we get

$$C'_+ = \cos \frac{\alpha}{2} e^{i(\beta+\gamma)/2} C_+ + i \sin \frac{\alpha}{2} e^{-i(\beta-\gamma)/2} C_-,$$

$$C'_- = i \sin \frac{\alpha}{2} e^{i(\beta-\gamma)/2} C_+ + \cos \frac{\alpha}{2} e^{-i(\beta+\gamma)/2} C_-. \quad (6.35)$$

So just starting from some assumptions about the properties of space, we have derived the amplitude transformation for any rotation at all. That means that if

† With a little work you can show that the frame x, y, z can also be brought into the frame x', y', z' by the following three rotations about the *original* axes: (1) rotate by the angle γ around the original z -axis; (2) rotate by the angle α around the original x -axis; (3) rotate by the angle β around the original z -axis.

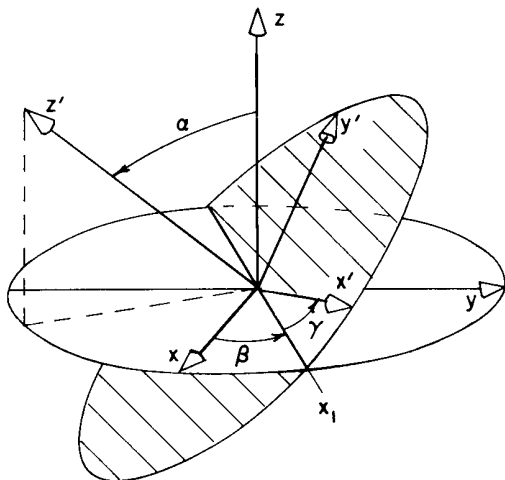


Fig. 6-9. The orientation of any coordinate frame x', y', z' relative to another frame x, y, z can be defined in terms of Euler's angles α, β, γ .

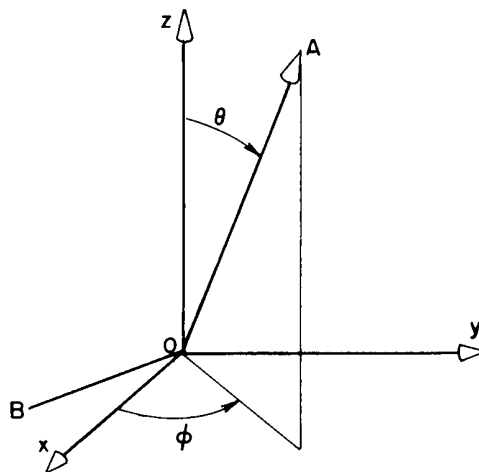


Fig. 6-10. An axis A defined by the polar angles θ and ϕ .

we know the amplitudes for any state of a spin one-half particle to go into the two beams of a Stern-Gerlach apparatus S , whose axes are x, y , and z , we can calculate what fraction would go into either beam of an apparatus T with the axes x', y' , and z' . In other words, if we have a state ψ of a spin one-half particle, whose amplitudes are $C_+ = \langle + | \psi \rangle$ and $C_- = \langle - | \psi \rangle$ to be “up” and “down” with respect to the z -axis of the x, y, z frame, we also know the amplitudes C'_+ and C'_- to be “up” and “down” with respect to the z' -axis of any other frame x', y', z' . The four coefficients in Eqs. (6.35) are the terms of the “transformation matrix” with which we can project the amplitudes of a spin one-half particle into any other coordinate system.

We will now work out a few examples to show you how it all works. Let's take the following simple question. We put a spin one-half atom through a Stern-Gerlach apparatus that transmits only the $(+z)$ state. What is the amplitude that it will be in the $(+x)$ state? The $+x$ axis is the same as the $+z'$ axis of a system rotated 90° about the y -axis. For this problem, then, it is simplest to use Eqs. (6.32)—although you could, of course, use the complete equations of (6.35). Since $C_+ = 1$ and $C_- = 0$, we get $C'_+ = 1/\sqrt{2}$. The probabilities are the absolute square of these amplitudes; there is a 50 percent chance that the particle will go through an apparatus that selects the $(+x)$ state. If we had asked about the $(-x)$ state the amplitude would have been $-1/\sqrt{2}$, which also gives a probability $1/2$ —as you would expect from the symmetry of space. So if a particle is in the $(+z)$ state, it is equally likely to be in $(+x)$ or $(-x)$, but with opposite phase.

There's no prejudice in y either. A particle in the $(+z)$ state has a 50-50 chance of being in $(+y)$ or in $(-y)$. However, for these (using the formula for rotating -90° about x), the amplitudes are $1/\sqrt{2}$ and $-i/\sqrt{2}$. In this case, the two amplitudes have a phase difference of 90° instead of 180° , as they did for the $(+x)$ and $(-x)$. In fact, that's how the distinction between x and y shows up.

As our final example, suppose that we know that a spin one-half particle is in a state ψ such that it is polarized “up” along some axis A , defined by the angles θ and ϕ in Fig. 6-10. We want to know the amplitude $\langle C_+ | \psi \rangle$ that the particle is “up” along z and the amplitude $\langle C_- | \psi \rangle$ that it is “down” along z . We can find these amplitudes by imagining that A is the z -axis of a system whose x -axis lies in some arbitrary direction—say in the plane formed by A and z . We can then bring the frame of A into x, y, z by three rotations. First, we make a rotation by $-\pi/2$ about the axis A , which brings the x -axis into the line B in the figure. Then we rotate by θ about line B (the new x -axis of frame A) to bring A to the z -axis. Finally, we rotate by the angle $(\pi/2 - \phi)$ about x . Remembering that we have only a $(+)$

state with respect to A , we get

$$C_+ = \cos \frac{\theta}{2} e^{-i\phi/2}, \quad C_- = \sin \frac{\theta}{2} e^{+i\phi/2}. \quad (6.36)$$

We would like, finally, to summarize the results of this chapter in a form that will be useful for our later work. First, we remind you that our primary result in Eqs. (6.35) can be written in another notation. Note that Eqs. (6.35) mean just the same thing as Eq. (6.4). That is, in Eqs. (6.35) the coefficients of $C_+ = \langle +S | \psi \rangle$ and $C_- = \langle -S | \psi \rangle$ are just the amplitudes $\langle jT | iS \rangle$ of Eq. (6.4)—the amplitudes that a particle in the i -state with respect to S will be in the j -state with respect to T (when the orientation of T with respect to S is given in terms of the angles α , β , and γ). We also called them R_{ji}^{TS} in Eq. (6.6). (We have a plethora of notations!) For example, $R_{-+}^{TS} = \langle -T | +S \rangle$ is the coefficient of C_+ in the formula for C_- , namely, $i \sin (\alpha/2) e^{i(\beta-\gamma)/2}$. We can, therefore, make a summary of our results in the form of a table, as we have done in Table 6-1.

It will occasionally be handy to have these amplitudes already worked out for some simple special cases. Let's let $R_z(\phi)$ stand for a rotation by the angle ϕ about the z -axis. We can also let it stand for the corresponding rotation matrix (omitting the subscripts i and j , which are to be implicitly understood). In the same spirit $R_x(\phi)$ and $R_y(\phi)$ will stand for rotations by the angle ϕ about the x -axis or the y -axis. We give in Table 6-2 the matrices—the tables of amplitudes $\langle jT | iS \rangle$ —which project the amplitudes from the S -frame into the T -frame, where T is obtained from S by the rotation specified.

Table 6-1

The amplitudes $\langle jT | iS \rangle$ for a rotation defined by the Euler angles α , β , γ of Fig. 6-9

$R_{ji}(\alpha, \beta, \gamma)$		
$\langle jT iS \rangle$	$+S$	$-S$
$+T$	$\cos \frac{\alpha}{2} e^{i(\beta+\gamma)/2}$	$i \sin \frac{\alpha}{2} e^{-i(\beta-\gamma)/2}$
$-T$	$i \sin \frac{\alpha}{2} e^{i(\beta-\gamma)/2}$	$\cos \frac{\alpha}{2} e^{-i(\beta+\gamma)/2}$

Table 6-2

The amplitudes $\langle jT | iS \rangle$ for a rotation $R(\phi)$ by the angle ϕ about the z -axis, x -axis, or y -axis

$R_z(\phi)$		
$\langle jT iS \rangle$	$+S$	$-S$
$+T$	$e^{i\phi/2}$	0
$-T$	0	$e^{-i\phi/2}$

$R_x(\phi)$		
$\langle jT iS \rangle$	$+S$	$-S$
$+T$	$\cos \phi/2$	$i \sin \phi/2$
$-T$	$i \sin \phi/2$	$\cos \phi/2$

$R_y(\phi)$		
$\langle jT iS \rangle$	$+S$	$-S$
$+T$	$\cos \phi/2$	$\sin \phi/2$
$-T$	$-\sin \phi/2$	$\cos \phi/2$

The Dependence of Amplitudes on Time

7-1 Atoms at rest; stationary states

We want now to talk a little bit about the behavior of probability amplitudes in time. We say a “little bit,” because the actual behavior in time necessarily involves the behavior in space as well. Thus, we get immediately into the most complicated possible situation if we are to do it correctly and in detail. We are always in the difficulty that we can either treat something in a logically rigorous but quite abstract way, or we can do something which is not at all rigorous but which gives us some idea of a real situation—postponing until later a more careful treatment. With regard to energy dependence, we are going to take the second course. We will make a number of statements. We will not try to be rigorous—but will just be telling you things that have been found out, to give you some feeling for the behavior of amplitudes as a function of time. As we go along, the precision of the description will increase, so don’t get nervous that we seem to be picking things out of the air. It is, of course, all out of the air—the air of experiment and of the imagination of people. But it would take us too long to go over the historical development, so we have to plunge in somewhere. We could plunge into the abstract and deduce everything—which you would not understand—or we could go through a large number of experiments to justify each statement. We choose to do something in between.

An electron alone in empty space can, under certain circumstances, have a certain definite energy. For example, if it is standing still (so it has no translational motion, no momentum, or kinetic energy), it has its rest energy. A more complicated object like an atom can also have a definite energy when standing still, but it could also be internally excited to another energy level. (We will describe later the machinery of this.) We can often think of an atom in an excited state as having a definite energy, but this is really only approximately true. An atom doesn’t stay excited forever because it manages to discharge its energy by its interaction with the electromagnetic field. So there is some amplitude that a new state is generated—with the atom in a lower state, and the electromagnetic field in a higher state, of excitation. The total energy of the system is the same before and after, but the energy of the *atom* is reduced. So it is not precise to say an excited atom has a *definite* energy; but it will often be convenient and not too wrong to say that it does.

[Incidentally, why does it go one way instead of the other way? Why does an atom radiate light? The answer has to do with entropy. When the energy is in the electromagnetic field, there are so many different ways it can be—so many different places where it can wander—that if we look for the equilibrium condition, we find that in the most probable situation the field is excited with a photon, and the atom is de-excited. It takes a very long time for the photon to come back and find that it can knock the atom back up again. It’s quite analogous to the classical problem: Why does an accelerating charge radiate? It isn’t that it “wants” to lose energy, because, in fact, when it radiates, the energy of the world is the same as it was before. Radiation or absorption goes in the direction of increasing *entropy*.]

Nuclei can also exist in different energy levels, and in an approximation which disregards the electromagnetic effects, we can say that a nucleus in an excited state stays there. Although we know that it doesn’t stay there forever, it is often useful to start out with an approximation which is somewhat idealized and easier to think about. Also it is often a legitimate approximation under certain circumstances. (When we first introduced the classical laws of a falling body, we did not include friction, but there is almost never a case in which there isn’t *some* friction.)

7-1 Atoms at rest; stationary states

7-2 Uniform motion

7-3 Potential energy; energy conservation

7-4 Forces; the classical limit

7-5 The “precession” of a spin one-half particle

Review: Chapter 17, Vol. I, *Space-Time*
Chapter 48, Vol. I, *Beats*

Then there are the subnuclear “strange particles,” which have various masses. But the heavier ones disintegrate into other light particles, so again it is not correct to say that they have a precisely definite energy. That would be true only if they lasted forever. So when we make the approximation that they have a definite energy, we are forgetting the fact that they must blow up. For the moment, then, we will intentionally forget about such processes and learn later how to take them into account.

Suppose we have an atom—or an electron, or any particle—which at rest would have a definite energy E_0 . By the energy E_0 we mean the mass of the whole thing times c^2 . This mass includes any internal energy; so an excited atom has a mass which is different from the mass of the same atom in the ground state. (The *ground state* means the state of lowest energy.) We will call E_0 the “energy at rest.”

For an atom *at rest*, the quantum mechanical *amplitude* to find an atom at a place is the *same everywhere*; it does *not* depend on position. This means, of course, that the *probability* of *finding* the atom anywhere is the same. But it means even more. The *probability* could be independent of position, and still the *phase* of the *amplitude* could vary from point to point. But for a particle at rest, the complete amplitude is identical everywhere. It does, however, depend on the *time*. For a particle in a state of definite energy E_0 , the amplitude to find the particle at (x, y, z) at the time t is

$$ae^{-i(E_0/\hbar)t}, \quad (7.1)$$

where a is some constant. The amplitude to be at any point in space is the same for all points, but depends on time according to (7.1). We shall simply assume this rule to be true.

Of course, we could also write (7.1) as

$$ae^{-i\omega t}, \quad (7.2)$$

with

$$\hbar\omega = E_0 = Mc^2,$$

where M is the rest mass of the atomic state, or particle. There are three different ways of specifying the energy: by the frequency of an amplitude, by the energy in the classical sense, or by the inertia. They are all equivalent; they are just different ways of saying the same thing.

You may be thinking that it is strange to think of a “particle” which has equal amplitudes to be found throughout all space. After all, we usually imagine a “particle” as a small object located “somewhere.” But don’t forget the uncertainty principle. If a particle has a definite energy, it has also a definite momentum. If the uncertainty in momentum is zero, the uncertainty relation, $\Delta p \Delta x = \hbar$, tells us that the uncertainty in the position must be infinite, and that is just what we are saying when we say that there is the same amplitude to find the particle at all points in space.

If the internal parts of an atom are in a different state with a different total energy, then the variation of the amplitude with time is different. If you don’t know in which state it is, there will be a certain amplitude to be in one state and a certain amplitude to be in another—and each of these amplitudes will have a different frequency. There will be an interference between these different components—like a beat-note—which can show up as a varying probability. Something will be “going on” inside of the atom—even though it is “at rest” in the sense that its center of mass is not drifting. However, if the atom has one definite energy, the amplitude is given by (7.1), and the absolute square of this amplitude does not depend on time. You see, then, that if a thing has a definite energy and if you ask any *probability* question about it, the answer is independent of time. Although the *amplitudes* vary with time, if the energy is *definite* they vary as an imaginary exponential, and the absolute value doesn’t change.

That’s why we often say that an atom in a definite energy level is in a *stationary state*. If you make any measurements of the things inside, you’ll find that nothing (in probability) will change in time. In order to have the probabilities change in

time, we have to have the interference of two amplitudes at two different frequencies, and that means that we cannot know what the energy is. The object will have one amplitude to be in a state of one energy and another amplitude to be in a state of another energy. That's the quantum mechanical description of something when its *behavior* depends on time.

If we have a "condition" which is a mixture of two different states with different energies, then the amplitude for each of the two states varies with time according to Eq. (7.2), for instance, as

$$e^{-i(E_1/\hbar)t} \quad \text{and} \quad e^{-i(E_2/\hbar)t}. \quad (7.3)$$

And if we have some combination of the two, we will have an interference. But notice that if we added a constant to both energies, it wouldn't make any difference. If somebody else were to use a different scale of energy in which all the energies were increased (or decreased) by a constant amount—say, by the amount A —then the amplitudes in the two states would, from his point of view, be

$$e^{-i(E_1+A)t/\hbar} \quad \text{and} \quad e^{-i(E_2+A)t/\hbar}. \quad (7.4)$$

All of his amplitudes would be multiplied by the same factor $e^{-i(A/\hbar)t}$, and all linear combinations, or interferences, would have the same factor. When we take the absolute squares to find the probabilities, all the answers would be the same. The choice of an origin for our energy scale makes no difference; we can measure energy from any zero we want. For relativistic purposes it is nice to measure the energy so that the rest mass is included, but for many purposes that aren't relativistic it is often nice to subtract some standard amount from all energies that appear. For instance, in the case of an atom, it is usually convenient to subtract the energy $M_s c^2$, where M_s is the mass of all the *separate* pieces—the nucleus and the electrons—which is, of course, different from the mass of the atom. For other problems it may be useful to subtract from all energies the amount $M_0 c^2$, where M_0 is the mass of the whole atom *in the ground state*; then the energy that appears is just the excitation energy of the atom. So, sometimes we may shift our zero of energy by some very large constant, but it doesn't make any difference, provided we shift all the energies in a particular calculation by the same constant. So much for a particle standing still.

7-2 Uniform motion

If we suppose that the relativity theory is right, a particle at rest in one inertial system can be in uniform motion in another inertial system. In the rest frame of the particle, the probability amplitude is the same for all x , y , and z but varies with t . The *magnitude* of the amplitude is the same for all t , but the *phase* depends on t . We can get a kind of a picture of the behavior of the amplitude if we plot lines of equal phase—say, lines of zero phase—as a function of x and t . For a particle at rest, these equal-phase lines are parallel to the x -axis and are equally spaced in the t -coordinate, as shown by the dashed lines in Fig. 7-1.

In a different frame— x' , y' , z' , t' —that is moving with respect to the particle in, say, the x -direction, the x' and t' coordinates of any particular point in space are related to x and t by the Lorentz transformation. This transformation can be represented graphically by drawing x' and t' axes, as is done in Fig. 7-1. (See Chapter 17, Vol. I, Fig. 17-2.) You can see that in the x' - t' system, points of equal phase† have a different spacing along the t' -axis, so the frequency of the time variation is different. Also there is a variation of the phase with x' , so the probability amplitude must be a function of x' .

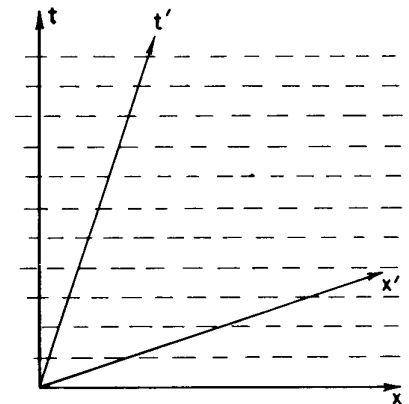


Fig. 7-1. Relativistic transformation of the amplitude of a particle at rest in the x - t systems.

† We are assuming that the phase should have the same value at corresponding points in the two systems. This is a subtle point, however, since the phase of a quantum mechanical amplitude is, to a large extent, arbitrary. A complete justification of this assumption requires a more detailed discussion involving interferences of two or more amplitudes.

Under a Lorentz transformation for the velocity v , say along the negative x -direction, the time t is related to the time t' by

$$t = \frac{t' - x'v/c^2}{\sqrt{1 - v^2/c^2}},$$

so our amplitude now varies as

$$e^{-(i/\hbar)E_0 t} = e^{-(i/\hbar)(E_0 t' / \sqrt{1 - v^2/c^2} - E_0 v x' / c^2 \sqrt{1 - v^2/c^2})}.$$

In the prime system it varies in space as well as in time. If we write the amplitude as

$$e^{-(i/\hbar)(E'_p t' - p' x')},$$

we see that $E'_p = E_0 / \sqrt{1 - v^2/c^2}$ is the energy computed classically for a particle of rest energy E_0 travelling at the velocity v , and $p' = E'_p v / c^2$ is the corresponding particle momentum.

You know that $x_\mu = (t, x, y, z)$ and $p_\mu = (E, p_x, p_y, p_z)$ are four-vectors, and that $p_\mu x_\mu = Et - \mathbf{p} \cdot \mathbf{x}$ is a scalar invariant. In the rest frame of the particle, $p_\mu x_\mu$ is just Et ; so if we transform to another frame, Et will be replaced by

$$E't' - \mathbf{p}' \cdot \mathbf{x}'.$$

Thus, the probability amplitude of a particle which has the momentum \mathbf{p} will be proportional to

$$e^{-(i/\hbar)(E_p t - \mathbf{p} \cdot \mathbf{x})}, \quad (7.5)$$

where E_p is the energy of the particle whose momentum is \mathbf{p} , that is,

$$E_p = \sqrt{(pc)^2 + E_0^2}, \quad (7.6)$$

where E_0 is, as before, the rest energy. For nonrelativistic problems, we can write

$$E_p = M_s c^2 + W_p, \quad (7.7)$$

where W_p is the energy over and above the rest energy $M_s c^2$ of the parts of the atom. In general, W_p would include both the kinetic energy of the atom as well as its binding or excitation energy, which we can call the “internal” energy. We would write

$$W_p = W_{\text{int}} + \frac{p^2}{2M}, \quad (7.8)$$

and the amplitudes would be

$$e^{-(i/\hbar)(W_p t - \mathbf{p} \cdot \mathbf{x})}. \quad (7.9)$$

Because we will generally be doing nonrelativistic calculations, we will use this form for the probability amplitudes.

Note that our relativistic transformation has given us the variation of the amplitude of an atom which moves in space without any additional assumptions. The wave number of the space variations is, from (7.9),

$$k = \frac{p}{\hbar}; \quad (7.10)$$

so the wavelength is

$$\lambda = \frac{2\pi}{k} = \frac{h}{p}. \quad (7.11)$$

This is the same wavelength we have used before for particles with the momentum p . This formula was first arrived at by de Broglie in just this way. For a moving particle, the *frequency* of the amplitude variations is still given by

$$\hbar\omega = W_p. \quad (7.12)$$

The absolute square of (7.9) is just 1, so for a particle in motion with a *definite energy*, the probability of finding it is the same everywhere and does not change with time. (It is important to notice that the amplitude is a *complex* wave. If we used a real sine wave, the square would vary from point to point, which would not be right.)

We know, of course, that there are situations in which particles move from place to place so that the probability depends on position and changes with time. How do we describe such situations? We can do that by considering amplitudes which are a superposition of two or more amplitudes for states of definite energy. We have already discussed this situation in Chapter 48 of Vol. I—even for probability amplitudes! We found that the sum of two amplitudes with different wave numbers k (that is, momenta) and frequencies ω (that is, energies) gives interference humps, or beats, so that the square of the amplitude varies with space and time. We also found that these beats move with the so-called “group velocity” given by

$$v_g = \frac{\Delta\omega}{\Delta k},$$

where Δk and $\Delta\omega$ are the differences between the wave numbers and frequencies for the two waves. For more complicated waves—made up of the sum of many amplitudes all near the same frequency—the group velocity is

$$v_g = \frac{d\omega}{dk}. \quad (7.13)$$

Taking $\omega = E_p/\hbar$ and $k = p/\hbar$, we see that

$$v_g = \frac{dE_p}{dp}. \quad (7.14)$$

Using Eq. (7.6), we have

$$\frac{dE_p}{dp} = c^2 \frac{p}{E_p}. \quad (7.15)$$

But $E_p = Mc^2$, so

$$\frac{dE_p}{dp} = \frac{p}{M}, \quad (7.16)$$

which is just the classical velocity of the particle. Alternatively, if we use the non-relativistic expressions, we have

$$\omega = \frac{W_p}{\hbar} \quad \text{and} \quad k = \frac{p}{\hbar},$$

and

$$\frac{d\omega}{dk} = \frac{dW}{dp} = \frac{d}{dp} \left(\frac{p^2}{2M} \right) = \frac{p}{M}, \quad (7.17)$$

which is again the classical velocity.

Our result, then, is that if we have several amplitudes for pure energy states of nearly the same energy, their interference gives “lumps” in the probability that move through space with a velocity equal to the velocity of a classical particle of that energy. We should remark, however, that when we say we can add two amplitudes of different wave number together to get a beat-note that will correspond to a moving particle, we have introduced something new—something that we cannot deduce from the theory of relativity. We said what the amplitude did for a particle standing still and then deduced what it would do if the particle were moving. But we *cannot* deduce from these arguments what would happen when there are *two* waves moving with different speeds. If we stop one, we cannot stop the other. So we have added tacitly the *extra* hypothesis that not only is (7.9) a *possible* solution, but that there can also be solutions with all kinds of p 's for the same system, and that the different terms will interfere.

7-3 Potential energy; energy conservation

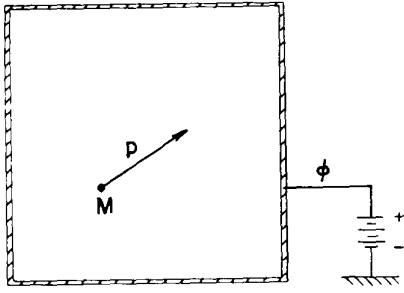


Fig. 7-2. A particle of mass M and momentum p in a region of constant potential.

Now we would like to discuss what happens when the energy of a particle can change. We begin by thinking of a particle which moves in a force field described by a potential. We discuss first the effect of a constant potential. Suppose that we have a large metal can which we have raised to some electrostatic potential ϕ , as in Fig. 7-2. If there are charged objects inside the can, their potential energy will be $q\phi$, which we will call V , and will be absolutely independent of position. Then there can be no change in the physics inside, because the constant potential doesn't make any difference so far as anything going on inside the can is concerned. Now there is no way we can deduce what the answer should be, so we must make a guess. The guess which works is more or less what you might expect: For the energy, we must use the sum of the potential energy V and the energy E_p —which is itself the sum of the internal and kinetic energies. The amplitude is proportional to

$$e^{-(i/\hbar)[(E_p + V)t - \mathbf{p} \cdot \mathbf{x}]} \quad (7.18)$$

The *general principle* is that the coefficient of t , which we may call ω , is always given by the *total energy* of the system: internal (or "mass") energy, plus kinetic energy, plus potential energy:

$$\hbar\omega = E_p + V. \quad (7.19)$$

Or, for nonrelativistic situations,

$$\hbar\omega = W_{\text{int}} + \frac{p^2}{2M} + V. \quad (7.20)$$

Now what about physical phenomena inside the box? If there are several different energy states, what will we get? The amplitude for each state has the same additional factor

$$e^{-(i/\hbar)Vt}$$

over what it would have with $V = 0$. That is just like a change in the zero of our energy scale. It produces an equal phase change in all amplitudes, but as we have seen before, this doesn't change any of the probabilities. All the physical phenomena are the same. (We have assumed that we are talking about different states of the same charged object, so that $q\phi$ is the same for all. If an object could change its charge in going from one state to another, we would have quite another result, but conservation of charge prevents this.)

So far, our assumption agrees with what we would expect for a change of energy reference level. But if it is really right, it should hold for a potential energy that is not just a constant. In general, V could vary in any arbitrary way with both time and space, and the complete result for the amplitude must be given in terms of a differential equation. We don't want to get concerned with the general case right now, but only want to get some idea about how some things happen, so we will think only of a potential that is constant in time and varies very slowly in space. Then we can make a comparison between the classical and quantum ideas.

Suppose we think of the situation in Fig. 7-3, which has two boxes held at the constant potentials ϕ_1 and ϕ_2 and a region in between where we will assume that the potential varies smoothly from one to the other. We imagine that some particle has an amplitude to be found in any one of the regions. We also assume that the momentum is large enough so that in any small region in which there are many wavelengths, the potential is nearly constant. We would then think that in any part of the space the amplitude ought to look like (7.18) with the appropriate V for that part of the space.

Let's think of a special case in which $\phi_1 = 0$, so that the potential energy there is zero, but in which $q\phi_2$ is negative, so that classically the particle would have more energy in the second box. Classically, it would be going faster in the second box—it would have more energy and, therefore, more momentum. Let's see how that might come out of quantum mechanics.

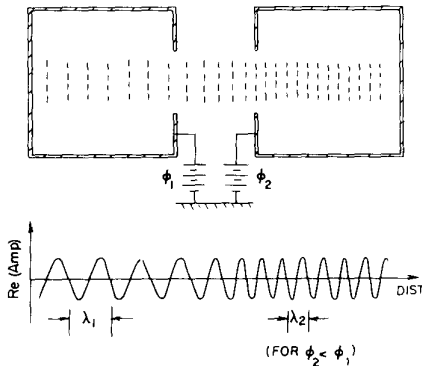


Fig. 7-3. The amplitude for a particle in transit from one potential to another.

With our assumption, the amplitude in the first box would be proportional to

$$e^{-(i/\hbar)[(W_{\text{int}} + p_1^2/2M + V_1)t - p_1 \cdot x]}, \quad (7.21)$$

and the amplitude in the second box would be proportional to

$$e^{-(i/\hbar)[(W_{\text{int}} + p_2^2/2M + V_2)t - p_2 \cdot x]}. \quad (7.22)$$

(Let's say that the internal energy is not being changed, but remains the same in both regions.) The question is: How do these two amplitudes match together through the region between the boxes?

We are going to suppose that the potentials are all constant in time—so that nothing in the conditions varies. We will then suppose that the variations of the amplitude (that is, its phase) have the same *frequency* everywhere—because, so to speak, there is nothing in the “medium” that depends on time. If nothing in the space is changing, we can consider that the wave in one region “generates” subsidiary waves all over space which will all oscillate at the same frequency—just as light waves going through materials at rest do not change their frequency. If the frequencies in (7.21) and (7.22) are the same, we must have that

$$W_{\text{int}} + \frac{p_1^2}{2M} + V_1 = W_{\text{int}} + \frac{p_2^2}{2M} + V_2. \quad (7.23)$$

Both sides are just the classical total energies, so Eq. (7.23) is a statement of the conservation of energy. In other words, the classical statement of the conservation of energy is equivalent to the quantum mechanical statement that the frequencies for a particle are everywhere the same if the conditions are not changing with time. It all fits with the idea that $\hbar\omega = E$.

In the special example that $V_1 = 0$ and V_2 is negative, Eq. (7.23) gives that p_2 is greater than p_1 , so the wavelength of the waves is shorter in region 2. The surfaces of equal phase are shown by the dashed lines in Fig. 7-3. We have also drawn a graph of the real part of the amplitude, which shows again how the wavelength decreases in going from region 1 to region 2. The group velocity of the waves, which is p/M , also increases in the way one would expect from the classical energy conservation, since it is just the same as Eq. (7.23).

There is an interesting special case where V_2 gets so large that $V_2 - V_1$ is greater than $p_1^2/2M$. Then p_2^2 , which is given by

$$p_2^2 = 2M \left[\frac{p_1^2}{2M} - V_2 + V_1 \right], \quad (7.24)$$

is *negative*. That means that p_2 is an imaginary number, say, ip' . Classically, we would say that the particle never gets into region 2—it doesn't have enough energy to climb the potential hill. Quantum mechanically, however, the amplitude is still given by Eq. (7.22); its space variation still goes as

$$e^{(i/\hbar)p_2 \cdot x}.$$

But if p_2 is imaginary, the space dependence becomes a real exponential. Say that the particle was initially going in the $+x$ -direction; then the amplitude would vary as

$$e^{-p'x/\hbar}. \quad (7.25)$$

The amplitude decreases rapidly with increasing x .

Imagine that the two regions at different potentials were very close together, so that the potential energy changed suddenly from V_1 to V_2 , as shown in Fig. 7-4(a). If we plot the real part of the probability amplitude, we get the dependence shown in part (b) of the figure. The wave in the first region corresponds to a particle trying to get into the second region, but the amplitude there falls off

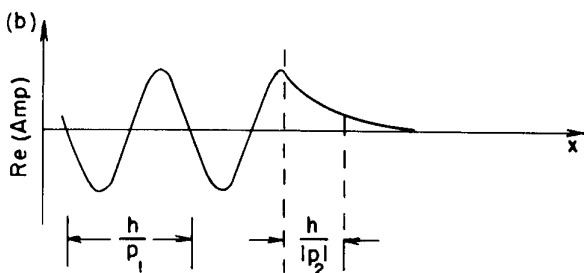
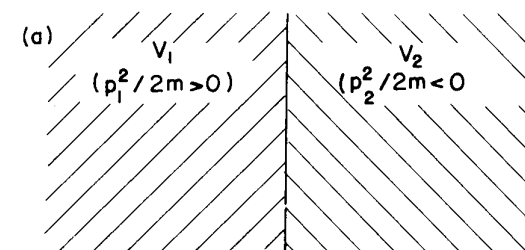


Fig. 7-4. The amplitude for a particle approaching a strongly repulsive potential.

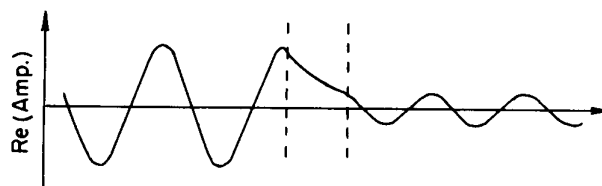
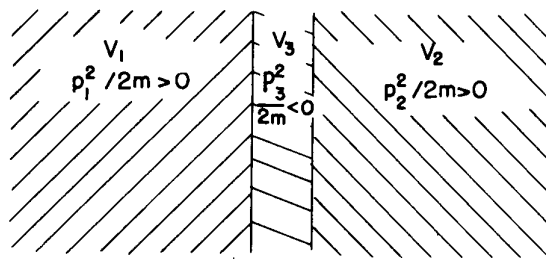


Fig. 7-5. The penetration of the amplitude through a potential barrier.

rapidly. There is some chance that it will be observed in the second region—where it could *never* get classically—but the amplitude is very small except right near the boundary. The situation is very much like what we found for the total internal reflection of light. The light doesn't normally get out, but we can observe it if we put something within a wavelength or two of the surface.

You will remember that if we put a second surface close to the boundary where light was totally reflected, we could get some light transmitted into the second piece of material. The corresponding thing happens to particles in quantum mechanics. If there is a narrow region with a potential V , so great that the classical kinetic energy would be negative, the particle would classically never get past. But quantum mechanically, the exponentially decaying amplitude can reach across the region and give a small probability that the particle will be found on the other side where the kinetic energy is again positive. The situation is illustrated in Fig. 7-5. This effect is called the quantum mechanical “penetration of a barrier.”

The barrier penetration by a quantum mechanical amplitude gives the explanation—or description—of the α -particle decay of a uranium nucleus. The potential energy of an α -particle, as a function of the distance from the center, is shown in Fig. 7-6(a). If one tried to shoot an α -particle with the energy E into the nucleus, it would feel an electrostatic repulsion from the nuclear charge z and would, classically, get no closer than the distance r_1 where its total energy is equal to the potential energy V . Closer in, however, the potential energy is much lower because of the strong attraction of the short-range nuclear forces. How is it then that in radioactive decay we find α -particles which started out inside the nucleus coming out with the energy E ? Because they start out with the energy E inside the nucleus and “leak” through the potential barrier. The probability amplitude is roughly as sketched in part (b) of Fig. 7-6, although actually the exponential decay is much larger than shown. It is, in fact, quite remarkable that the mean life of an α -particle in the uranium nucleus is as long as $4\frac{1}{2}$ billion years, when the natural oscillations inside the nucleus are so extremely rapid—about 10^{22} per sec! How can one get a number like 10^9 years from 10^{-22} sec? The answer is that the exponential gives the tremendously small factor of about e^{-45} —which gives the very small, though definite, probability of leakage. Once the α -particle is in the nucleus, there is almost no amplitude at all for finding it outside; however, if you take many nuclei and wait long enough, you may be lucky and find one that has come out.

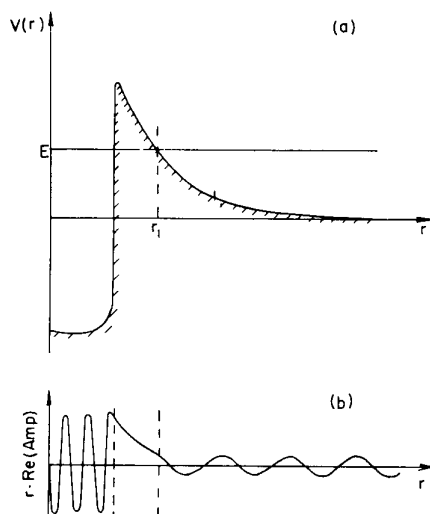


Fig. 7-6. (a) The potential function for an α -particle in a uranium nucleus. (b) The qualitative form of the probability amplitude.

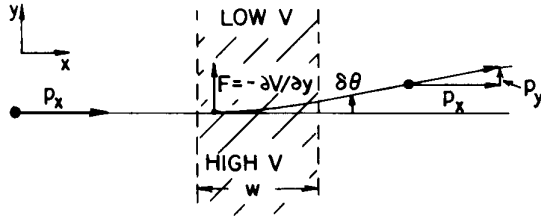


Fig. 7-7. The deflection of a particle by a transverse potential gradient.

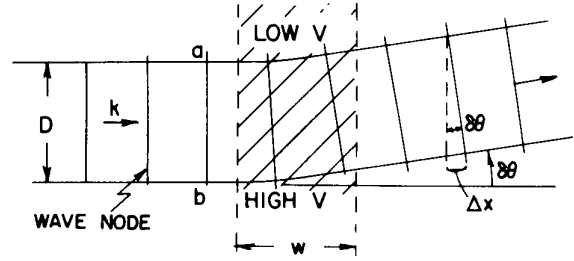


Fig. 7-8. The probability amplitude in a region with a transverse potential gradient.

7-4 Forces; the classical limit

Suppose that we have a particle moving along and passing through a region where there is a potential that varies at right angles to the motion. Classically, we would describe the situation as sketched in Fig. 7-7. If the particle is moving along the x -direction and enters a region where there is a potential that varies with y , the particle will get a transverse acceleration from the force $F = -\partial V/\partial y$. If the force is present only in a limited region of width w , the force will act only for the time w/v . The particle will be given the transverse momentum

$$p_y = F \frac{w}{v}.$$

The angle of deflection $\delta\theta$ is then

$$\delta\theta = \frac{p_y}{p} = \frac{Fw}{pv},$$

where p is the initial momentum. Using $-\partial V/\partial y$ for F , we get

$$\delta\theta = -\frac{w}{pv} \frac{\partial V}{\partial y}. \quad (7.26)$$

It is now up to us to see if our idea that the waves go as (7.20) will explain the same result. We look at the same thing quantum mechanically, assuming that everything is on a very large scale compared with a wavelength of our probability amplitudes. In any small region we can say that the amplitude varies as

$$e^{-(i/\hbar)[(W+p^2/2M+V)t-\mathbf{p}\cdot\mathbf{x}]}. \quad (7.27)$$

Can we see that this will also give rise to a deflection of the particle when V has a transverse gradient? We have sketched in Fig. 7-8 what the waves of probability amplitude will look like. We have drawn a set of "wave nodes" which you can think of as surfaces where the phase of the amplitude is zero. In every small region, the wavelength—the distance between successive nodes—is

$$\lambda = \frac{h}{p},$$

where p is related to V through

$$W + \frac{p^2}{2M} + V = \text{const.} \quad (7.28)$$

In the region where V is larger, p is smaller, and the wavelength is longer. So the angle of the wave nodes gets changed as shown in the figure.

To find the change in angle of the wave nodes we notice that for the two paths a and b in Fig. 7-8 there is a difference of potential $\Delta V = (\partial V/\partial y)D$, so there is a difference Δp in the momentum along the two tracks which can be

obtained from (7.28):

$$\Delta \left(\frac{p^2}{2M} \right) = \frac{p}{M} \Delta p = -\Delta V. \quad (7.29)$$

The wave number p/\hbar is, therefore, different along the two paths, which means that the phase is advancing at a different rate. The difference in the rate of increase of phase is $\Delta k = \Delta p/\hbar$, so the accumulated phase difference in the total distance w is

$$\Delta(\text{phase}) = \Delta k \cdot w = \frac{\Delta p}{\hbar} \cdot w = -\frac{M}{p\hbar} \Delta V \cdot w. \quad (7.30)$$

This is the amount by which the phase on path b is “ahead” of the phase on path a as the wave leaves the strip. But outside the strip, a phase advance of this amount corresponds to the wave node being ahead by the amount

$$\Delta x = \frac{\lambda}{2\pi} \Delta(\text{phase}) = \frac{\hbar}{p} \Delta(\text{phase})$$

or

$$\Delta x = -\frac{M}{p^2} \Delta V \cdot w. \quad (7.31)$$

Referring to Fig. 7-8, we see that the new wavefronts will be at the angle $\delta\theta$ given by

$$\Delta x = D \delta\theta; \quad (7.32)$$

so we have

$$D \delta\theta = -\frac{M}{p^2} \Delta V \cdot w. \quad (7.33)$$

This is identical to Eq. (7.26) if we replace p/m by v and $\Delta V/D$ by $\partial V/\partial y$.

The result we have just got is correct only if the potential variations are slow and smooth—in what we call the *classical limit*. We have shown that under these conditions we will get the same particle motions we get from $F = ma$, provided we assume that a potential contributes a phase to the probability amplitude equal to Vt/\hbar . *In the classical limit, the quantum mechanics will agree with Newtonian mechanics.*

7-5 The “precession” of a spin one-half particle

Notice that we have not assumed anything special about the potential energy—it is just that energy whose derivative gives a force. For instance, in the Stern-Gerlach experiment we had the energy $U = -\boldsymbol{\mu} \cdot \mathbf{B}$, which gives a force if \mathbf{B} has a spatial variation. If we wanted to give a quantum mechanical description, we would have said that the particles in one beam had an energy that varied one way and that those in the other beam had an opposite energy variation. (We could put the magnetic energy U into the potential energy V or into the “internal” energy W ; it doesn’t matter.) Because of the energy variation, the waves are refracted, and the beams are bent up or down. (We see now that quantum mechanics would give us the same bending as we would compute from the classical mechanics.)

From the dependence of the amplitude on potential energy we would also expect that if a particle sits in a uniform magnetic field along the z -direction, its probability amplitude must be changing with time according to

$$e^{-(i/\hbar)(-\mu_z B)t}.$$

(We can consider that this is, in effect, a definition of μ_z .) In other words, if we place a particle in a uniform field B for a time τ , its probability amplitude will be multiplied by

$$e^{-(i/\hbar)(-\mu_z B)\tau}$$

over what it would be in no field. Since for a spin one-half particle, μ_z can be either plus or minus some number, say μ , the two possible states in a uniform field would have their phases changing at the same rate but in opposite directions. The two amplitudes get multiplied by

$$e^{\pm(i/\hbar)\mu B\tau}. \quad (7.34)$$

This result has some interesting consequences. Suppose we have a spin one-half particle in some state that is not purely spin up or spin down. We can describe its condition in terms of the amplitudes to be in the pure up and pure down states. But in a magnetic field, these two states will have phases changing at a different rate. So if we ask some question about the amplitudes, the answer will depend on how long it has been in the field.

As an example, we consider the disintegration of the muon in a magnetic field. When muons are produced as disintegration products of π -mesons, they are polarized (in other words, they have a preferred spin direction). The muons, in turn, disintegrate—in about 2.2 microseconds on the average—emitting an electron and two neutrinos:

$$\mu \rightarrow e + \nu + \bar{\nu}.$$

In this disintegration it turns out that (for at least the highest energies) the electrons are emitted preferentially in the direction opposite to the spin direction of the muon.

Suppose then that we consider the experimental arrangement shown in Fig. 7-9. If polarized muons enter from the left and are brought to rest in a block of material at A , they will, a little while later, disintegrate. The electrons emitted will, in general, go off in all possible directions. Suppose, however, that the muons all enter the stopping block at A with their spins in the x -direction. Without a magnetic field there would be some angular distribution of decay directions; we would like to know how this distribution is changed by the presence of the magnetic field. We expect that it may vary in some way with time. We can find out what happens by asking, for any moment, what the amplitude is that the muon will be found in the $(+x)$ state.

We can state the problem in the following way: A muon is known to have its spin in the $+x$ -direction at $t = 0$; what is the amplitude that it will be in the same state at the time τ ? Now we do not have any rule for the behavior of a spin one-half particle in a magnetic field at right angles to the spin, but we do know what happens to the spin up and spin down states with respect to the field—their amplitudes get multiplied by the factor (7.34). Our procedure then is to choose the representation in which the base states are spin up and spin down with respect to the z -direction (the field direction). Any question can then be expressed with reference to the amplitudes for these states.

Let's say that $\psi(t)$ represents the muon state. When it enters the block A , its state is $\psi(0)$, and we want to know $\psi(\tau)$ at the later time τ . If we represent the two base states by $(+z)$ and $(-z)$ we know the two amplitudes $\langle +z | \psi(0) \rangle$ and $\langle -z | \psi(0) \rangle$ —we know these amplitudes because we know that $\psi(0)$ represents a state with the spin in the $(+x)$ state. From the results of the last chapter, these amplitudes are†

$$\langle +z | +x \rangle = C_+ = \frac{1}{\sqrt{2}} \quad (7.35)$$

and

$$\langle -z | +x \rangle = C_- = \frac{1}{\sqrt{2}}.$$

They happen to be equal. Since these amplitudes refer to the condition at $t = 0$, let's call them $C_+(0)$ and $C_-(0)$.

† If you skipped Chapter 6, you can just take (7.35) as an underived rule for now. We will give later (in Chapter 10) a more complete discussion of spin precession, including a derivation of these amplitudes.

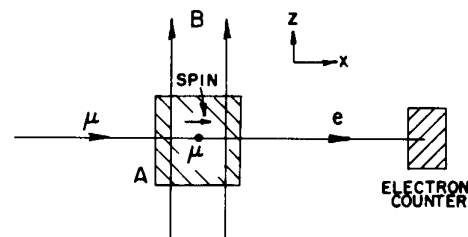


Fig. 7-9. A muon-decay experiment.

Now we know what happens to these two amplitudes with time. Using (7.34), we have

$$C_+(t) = C_+(0)e^{-(i/\hbar)\mu B t} \quad (7.36)$$

and

$$C_-(t) = C_-(0)e^{+(i/\hbar)\mu B t}$$

But if we know $C_+(t)$ and $C_-(t)$, we have all there is to know about the condition at t . The only trouble is that what we *want* to know is the probability that at t the spin will be in the $+x$ -direction. Our general rules can, however, take care of this problem. We write that the amplitude to be in the $(+x)$ state at time t , which we may call $A_+(t)$, is

$$A_+(t) = \langle +x | \psi(t) \rangle = \langle +x | +z \rangle \langle +z | \psi(t) \rangle + \langle +x | -z \rangle \langle -z | \psi(t) \rangle$$

or

$$A_+(t) = \langle +x | +z \rangle C_+(t) + \langle +x | -z \rangle C_-(t). \quad (7.37)$$

Again using the results of the last chapter—or better the equality $\langle \phi | x \rangle = \langle x | \phi \rangle^*$ from Chapter 5—we know that

$$\langle +x | +z \rangle = \frac{1}{\sqrt{2}} \quad \langle +x | -z \rangle = \frac{1}{\sqrt{2}}.$$

So we know all the quantities in Eq. (7.37). We get

$$A_+(t) = \frac{1}{2}e^{(i/\hbar)\mu B t} + \frac{1}{2}e^{-(i/\hbar)\mu B t},$$

or

$$A_+(t) = \cos \frac{\mu B}{\hbar} t.$$

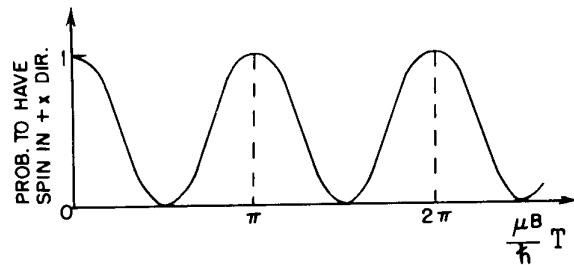
A particularly simple result! Notice that the answer agrees with what we expect for $t = 0$. We get $A_+(0) = 1$, which is right, because we assumed that the muon was in the $(+x)$ state at $t = 0$.

The probability P_+ that the muon will be found in the $(+x)$ state at t is $(A_+)^2$ or

$$P_+ = \cos^2 \frac{\mu B t}{\hbar}.$$

The probability oscillates between zero and one, as shown in Fig. 7-10. Note that the probability returns to one for $\mu B t / \hbar = \pi$ (not 2π). Because we have squared the cosine function, the probability repeats itself with the frequency $2\mu B / \hbar$.

Fig. 7-10. Time dependence of the probability that a spin one-half particle will be in a $(+)$ state with respect to the x -axis.



Thus, we find that the chance of catching the decay electron in the electron counter of Fig. 7-9 varies periodically with the length of time the muon has been sitting in the magnetic field. The frequency depends on the magnetic moment μ . The magnetic moment of the muon has, in fact, been measured in just this way.

We can, of course, use the same method to answer any other questions about the muon decay. For example, how does the chance of detecting a decay electron

in the y -direction at 90° to the x -direction but still at right angles to the field depend on t ? If you work it out, the amplitude to be in the $(+y)$ state varies as $\cos^2 \{(\mu B t / \hbar) - \pi/4\}$, which oscillates with the same period but reaches its maximum one-quarter cycle later, when $\mu B t / \hbar = \pi/4$. In fact, what is happening is that as time goes on, the muon goes through a succession of states which correspond to complete polarization in a direction that is continually rotating about the z -axis. We can describe this by saying that *the spin is precessing* at the frequency

$$\omega_p = \frac{2\mu B}{\hbar} . \quad (7.38)$$

You can begin to see the form that our quantum mechanical description will take when we are describing how things behave in time.

The Hamiltonian Matrix

8-1 Amplitudes and vectors

Before we begin the main topic of this chapter, we would like to describe a number of mathematical ideas that are used a lot in the literature of quantum mechanics. Knowing them will make it easier for you to read other books or papers on the subject. The first idea is the close mathematical resemblance between the equations of quantum mechanics and those of the scalar product of two vectors. You remember that if χ and ϕ are two states, the amplitude to start in ϕ and end up in χ can be written as a sum over a complete set of base states of the amplitude to go from ϕ into one of the base states and then from that base state out again into χ :

$$\langle \chi | \phi \rangle = \sum_{\text{all } i} \langle \chi | i \rangle \langle i | \phi \rangle. \quad (8.1)$$

We explained this in terms of a Stern-Gerlach apparatus, but we remind you that there is no need to have the apparatus. Equation (8.1) is a mathematical law that is just as true whether we put the filtering equipment in or not—it is not always necessary to imagine that the apparatus is there. We can think of it simply as a formula for the amplitude $\langle \chi | \phi \rangle$.

We would like to compare Eq. (8.1) to the formula for the dot product of two vectors \mathbf{B} and \mathbf{A} . If \mathbf{B} and \mathbf{A} are ordinary vectors in three dimensions, we can write the dot product this way:

$$\sum_{\text{all } i} (\mathbf{B} \cdot \mathbf{e}_i)(\mathbf{e}_i \cdot \mathbf{A}), \quad (8.2)$$

with the understanding that the symbol \mathbf{e}_i stands for the three unit vectors in the x , y , and z -directions. Then $\mathbf{B} \cdot \mathbf{e}_1$ is what we ordinarily call B_x ; $\mathbf{B} \cdot \mathbf{e}_2$ is what we ordinarily call B_y ; and so on. So Eq. (8.2) is equivalent to

$$B_x A_x + B_y A_y + B_z A_z,$$

which is the dot product $\mathbf{B} \cdot \mathbf{A}$.

Comparing Eqs. (8.1) and (8.2), we can see the following analogy: The states χ and ϕ correspond to the two vectors \mathbf{A} and \mathbf{B} . The base states i correspond to the special vectors \mathbf{e}_i to which we refer all other vectors. Any vector can be represented as a linear combination of the three “base vectors” \mathbf{e}_i . Furthermore, if you know the coefficients of each “base vector” in this combination—that is, its three components—you know everything about a vector. In a similar way, any quantum mechanical state can be described completely by the amplitude $\langle i | \phi \rangle$ to go into the base states; and if you know these coefficients, you know everything there is to know about the state. Because of this close analogy, what we have called a “state” is often also called a “state vector.”

Since the base vectors \mathbf{e}_i are all at right angles, we have the relation

$$\mathbf{e}_i \cdot \mathbf{e}_j = \delta_{ij}. \quad (8.3)$$

This corresponds to the relations (5.25) among the base states i ,

$$\langle i | j \rangle = \delta_{ij}. \quad (8.4)$$

You see now why one says that the base states i are all “orthogonal.”

8-1 Amplitudes and vectors

8-2 Resolving state vectors

8-3 What are the base states of the world?

8-4 How states change with time

8-5 The Hamiltonian matrix

8-6 The ammonia molecule

Review: Chapter 49, Vol. I, Modes

There is one minor difference between Eq. (8.1) and the dot product. We have that

$$\langle \phi | \chi \rangle = \langle \chi | \phi \rangle^*. \quad (8.5)$$

But in vector algebra,

$$\mathbf{A} \cdot \mathbf{B} = \mathbf{B} \cdot \mathbf{A}.$$

With the complex numbers of quantum mechanics we have to keep straight the order of the terms, whereas in the dot product, the order doesn't matter.

Now consider the following vector equation:

$$\mathbf{A} = \sum_i \mathbf{e}_i (\mathbf{e}_i \cdot \mathbf{A}). \quad (8.6)$$

It's a little unusual, but correct. It means the same thing as

$$\mathbf{A} = \sum_i A_i \mathbf{e}_i = A_x \mathbf{e}_x + A_y \mathbf{e}_y + A_z \mathbf{e}_z. \quad (8.7)$$

Notice, though, that Eq. (8.6) involves a quantity which is *different* from a dot product. A dot product is just a *number*, whereas Eq. (8.6) is a *vector* equation. One of the great tricks of vector analysis was to abstract away from the equations the idea of a *vector* itself. One might be similarly inclined to abstract a thing that is the analog of a "vector" from the quantum mechanical formula Eq. (8.1)—and one can indeed. We remove the $\langle \chi |$ from both sides Eq. (8.1) and write the following equation (don't get frightened—it's just a notation and in a few minutes you will find out what the symbols mean):

$$| \phi \rangle = \sum_i | i \rangle \langle i | \phi \rangle. \quad (8.8)$$

One thinks of the bracket $\langle \chi | \phi \rangle$ as being divided into two pieces. The second piece $| \phi \rangle$ is often called a *ket*, and the first piece $\langle \chi |$ is called a *bra* (put together, they make a "bra-ket"—a notation proposed by Dirac); the half-symbols $\langle \chi |$ and $| \phi \rangle$ are also called *state vectors*. In any case, they are *not* numbers, and, in general, we want the results of our calculations to come out as numbers; so such "unfinished" quantities are only part-way steps in our calculations.

It happens that until now we have written all our results in terms of numbers. How have we managed to avoid vectors? It is amusing to note that even in ordinary vector algebra we *could* make all equations involve only numbers. For instance, instead of a vector equation like

$$\mathbf{F} = m\mathbf{a},$$

we could always have written

$$\mathbf{C} \cdot \mathbf{F} = \mathbf{C} \cdot (m\mathbf{a}).$$

We have then an equation between dot products that is true for *any* vector \mathbf{C} . But if it is true for any \mathbf{C} , it hardly makes sense at all to keep writing the \mathbf{C} !

Now look at Eq. (8.1). It is an equation that is true for *any* χ . So to save writing, we should just leave *out* the χ and write Eq. (8.8) instead. It has the same information *provided* we understand that it should always be "finished" by "multiplying on the left by"—which simply means reinserting—some $\langle \chi |$ on both sides. So Eq. (8.8) means exactly the same thing as Eq. (8.1)—no more, no less. When you want numbers, you put in the $\langle \chi |$ you want.

Maybe you have already wondered about the ϕ in Eq. (8.8). Since the equation is true for *any* ϕ , why do we keep *it*? Indeed, Dirac suggests that the ϕ also can just as well be abstracted away, so that we have only

$$| = \sum_i | i \rangle \langle i |. \quad (8.9)$$

And this is the great law of quantum mechanics! (There is no analog in vector analysis.) It says that if you put *in* any two states χ and ϕ on the left and right of both sides, you *get back* Eq. (8.1). It is not really very useful, but it's a nice reminder that the equation is true for any two states.

8-2 Resolving state vectors

Let's look at Eq. (8.8) again; we can think of it in the following way. Any state vector $|\phi\rangle$ can be represented as a linear combination with suitable coefficients of a set of base "vectors"—or, if you prefer, as a superposition of "unit vectors" in suitable proportions. To emphasize that the coefficients $\langle i|\phi\rangle$ are just ordinary (complex) numbers, suppose we write

$$\langle i|\phi\rangle = C_i.$$

Then Eq. (8.8) is the same as

$$|\phi\rangle = \sum_i |i\rangle C_i. \quad (8.10)$$

We can write a similar equation for any other state vector, say $|\chi\rangle$, with, of course, different coefficients—say D_i . Then we have

$$|\chi\rangle = \sum_i |i\rangle D_i. \quad (8.11)$$

The D_i are just the amplitudes $\langle i|\chi\rangle$.

Suppose we had started by abstracting the ϕ from Eq. (8.1). We would have had

$$\langle\chi| = \sum_i \langle\chi|i\rangle\langle i|. \quad (8.12)$$

Remembering that $\langle\chi|i\rangle = \langle i|\chi\rangle^*$, we can write this as

$$\langle\chi| = \sum_i D_i^* \langle i|. \quad (8.13)$$

Now the interesting thing is that we can just *multiply* Eq. (8.13) and Eq. (8.10) to get back $\langle\chi|\phi\rangle$. When we do that, we have to be careful of the summation indices, because they are quite distinct in the two equations. Let's first rewrite Eq. (8.13) as

$$\langle\chi| = \sum_j D_j^* \langle j|,$$

which changes nothing. Then putting it together with Eq. (8.10), we have

$$\langle\chi|\phi\rangle = \sum_{ij} D_j^* \langle j|i\rangle C_i. \quad (8.14)$$

Remember, though, that $\langle j|i\rangle = \delta_{ij}$, so that in the sum we have left only the terms with $j = i$. We get

$$\langle\chi|\phi\rangle = \sum_i D_i^* C_i, \quad (8.15)$$

where, of course, $D_i^* = \langle i|\chi\rangle^* = \langle\chi|i\rangle$, and $C_i = \langle i|\phi\rangle$. Again we see the close analogy with the dot product

$$\mathbf{A} \cdot \mathbf{B} = \sum_i A_i B_i.$$

The only difference is the complex conjugate on D_i . So Eq. (8.15) says that if the state vectors $\langle\chi|$ and $|\phi\rangle$ are expanded in terms of the base vectors $\langle i|$ or $|i\rangle$, the amplitude to go from ϕ to χ is given by the kind of dot product in Eq. (8.15). This equation is, of course, just Eq. (8.1) written with different symbols. So we have just gone in a circle to get used to the new symbols.

We should perhaps emphasize again that while space vectors in three dimensions are described in terms of *three* orthogonal unit vectors, the base vectors $|i\rangle$ of the quantum mechanical states must range over the complete set applicable to any particular problem. Depending on the situation, two, or three, or five, or an infinite number of base states may be involved.

We have also talked about what happens when particles go through an apparatus. If we start the particles out in a certain state ϕ , then send them through

an apparatus, and afterward make a measurement to see if they are in state χ , the result is described by the amplitude

$$\langle \chi | A | \phi \rangle. \quad (8.16)$$

Such a symbol doesn't have a close analog in vector algebra. (It is closer to tensor algebra, but the analogy is not particularly useful.) We saw in Chapter 5, Eq. (5.32), that we could write (8.16) as

$$\langle \chi | A | \phi \rangle = \sum_{ij} \langle \chi | i \rangle \langle i | A | j \rangle \langle j | \phi \rangle. \quad (8.17)$$

This is just an example of the fundamental rule Eq. (8.9), used twice.

We also found that if another apparatus B was added in series with A , then we could write

$$\langle \chi | BA | \phi \rangle = \sum_{ijk} \langle \chi | i \rangle \langle i | B | j \rangle \langle j | A | k \rangle \langle k | \phi \rangle. \quad (8.18)$$

Again, this comes directly from Dirac's method of writing Eq. (8.9)—remember that we can always place a bar ($|$), which is just like the factor 1, between B and A .

Incidentally, we can think of Eq. (8.17) in another way. Suppose we think of the particle entering apparatus A in the state ϕ and coming out of A in the state ψ ("psi"). In other words, we could ask ourselves this question: Can we find a ψ such that the amplitude to get from ψ to χ is always identically and everywhere the same as the amplitude $\langle \chi | A | \phi \rangle$? The answer is yes. We want Eq. (8.17) to be replaced by

$$\langle \chi | \psi \rangle = \sum_i \langle \chi | i \rangle \langle i | \psi \rangle. \quad (8.19)$$

We can clearly do this if

$$\langle i | \psi \rangle = \sum_j \langle i | A | j \rangle \langle j | \phi \rangle = \langle i | A | \phi \rangle, \quad (8.20)$$

which determines ψ . "But it doesn't determine ψ ," you say; "it only determines $\langle i | \psi \rangle$." However, $\langle i | \psi \rangle$ *does* determine ψ , because if you have all the coefficients that relate ψ to the base states i , then ψ is uniquely defined. In fact, we can play with our notation and write the last term of Eq. (8.20) as

$$\langle i | \psi \rangle = \sum_j \langle i | j \rangle \langle j | A | \phi \rangle. \quad (8.21)$$

Then, since this equation is true for all i , we can write simply

$$| \psi \rangle = \sum_j | j \rangle \langle j | A | \phi \rangle. \quad (8.22)$$

Then we can say: "The state ψ is what we get if we start with ϕ and go through the apparatus A ."

One final example of the tricks of the trade. We start again with Eq. (8.17). Since it is true for any χ and ϕ , we can drop them both! We then get†

$$A = \sum_{ij} | i \rangle \langle i | A | j \rangle \langle j |. \quad (8.23)$$

What does it mean? It means no more, no less, than what you get if you put back the ϕ and χ . As it stands, it is an "open" equation and incomplete. If we multiply it "on the left" by $| \phi \rangle$, it becomes

$$A | \phi \rangle = \sum_{ij} | i \rangle \langle i | A | j \rangle \langle j | \phi \rangle, \quad (8.24)$$

† You might think we should write $|A|$ instead of just A . But then it would look like the symbol for "absolute value of A ," so the bars are usually dropped. In general, the bar ($|$) behaves much like the factor one.

which is just Eq. (8.22) all over again. In fact, we could have just dropped the j 's from that equation and written

$$|\psi\rangle = A|\phi\rangle. \quad (8.25)$$

The symbol A is neither an amplitude, nor a vector; it is a new kind of thing called an *operator*. It is something which “operates on” a state to produce a new state—Eq. (8.25) says that $|\psi\rangle$ is what results if A operates on $|\phi\rangle$. Again, it is still an open equation until it is completed with some bra like $\langle x|$ to give

$$\langle x|\psi\rangle = \langle x|A|\phi\rangle. \quad (8.26)$$

The operator A is, of course, described completely if we give the matrix of amplitudes $\langle i|A|j\rangle$ —also written A_{ij} —in terms of any set of base vectors.

We have really added nothing new with all of this new mathematical notation. One reason for bringing it all up was to show you the way of writing pieces of equations, because in many books you will find the equations written in the incomplete forms, and there's no reason for you to be paralyzed when you come across them. If you prefer, you can always add the missing pieces to make an equation between numbers that will look like something more familiar.

Also, as you will see, the “bra” and “ket” notation is a very convenient one. For one thing, we can from now on identify a state by giving its state vector. When we want to refer to a state of definite momentum p we can say: “the state $|p\rangle$ ”. Or we may speak of some arbitrary state $|\psi\rangle$. For consistency we will always use the ket, writing $|\psi\rangle$, to identify a state. (It is, of course an arbitrary choice; we could equally well have chosen to use the bra, $\langle\psi|$.)

8-3 What are the base states of the world?

We have discovered that any state in the world can be represented as a superposition—a linear combination with suitable coefficients—of base states. You may ask, first of all, *what* base states? Well, there are many different possibilities. You can, for instance, project a spin in the z -direction or in some other direction. There are many, many different *representations*, which are the analogs of the different *coordinate systems* one can use to represent ordinary vectors. Next, *what* coefficients? Well, that depends on the physical circumstances. Different sets of coefficients correspond to different physical conditions. The important thing to know about is the “space” in which you are working—in other words, what the base states mean physically. So the first thing you have to know about, in general, is what the base states are like. Then you can understand how to describe a situation in terms of these base states.

We would like to look ahead a little and speak a bit about what the general quantum mechanical description of nature is going to be—in terms of the now current ideas of physics, anyway. First, one decides on a particular representation for the base states—different representations are always possible. For example, for a spin one-half particle we can use the plus and minus states with respect to the z -axis. But there's nothing special about the z -axis—you can take any other axis you like. For consistency we'll always pick the z -axis, however. Suppose we begin with a situation with one electron. In addition to the two possibilities for the spin (“up” and “down” along the z -direction), there is also the momentum of the electron. We pick a set of base states, each corresponding to one value of the momentum. What if the electron doesn't have a definite momentum? That's all right; we're just saying what the *base* states are. If the electron hasn't got a definite momentum, it has some amplitude to have one momentum and another amplitude to have another momentum, and so on. And if it is not necessarily spinning up, it has some amplitude to be spinning up going at this momentum, and some amplitude to be spinning down going at that momentum, and so on. The complete description of an electron, *so far as we know*, requires only that the base states be described by the *momentum* and the *spin*. So one acceptable set of base states $|i\rangle$ for a single electron refer to different values of the momentum and

whether the spin is up or down. Different mixtures of amplitudes—that is, different combinations of the C 's describe different circumstances. What any particular electron is doing is described by telling with what amplitude it has an up-spin or a down-spin and one momentum or another—for all possible momenta. So you can see what is involved in a complete quantum mechanical description of a single electron.

What about systems with more than one electron? Then the base states get more complicated. Let's suppose that we have two electrons. We have, first of all, four possible states with respect to spin: both electrons spinning up, the first one down and the second one up, the first one up and the second one down, or both down. Also we have to specify that the first electron has the momentum p_1 , and the second electron, the momentum p_2 . The base states for two electrons require the specification of two momenta and two spin characters. With seven electrons, we have to specify seven of each.

If we have a proton and an electron, we have to specify the spin direction of the proton and its momentum, and the spin direction of the electron and its momentum. At least that's approximately true. *We do not really know* what the correct representation is for the world. It is all very well to start out by supposing that if you specify the spin in the electron and its momentum, and likewise for a proton, you will have the base states; but what about the "guts" of the proton? Let's look at it this way. In a hydrogen atom which has one proton and one electron, we have many different base states to describe—up and down spins of the proton and electron and the various possible momenta of the proton and electron. Then there are different combinations of amplitudes C_i which together describe the character of the hydrogen atom in different states. But suppose we look at the whole hydrogen atom as a "particle." If we didn't know that the hydrogen atom was made out of a proton and an electron, we might have started out and said: "Oh, I know what the base states are—they correspond to a particular momentum of the hydrogen atom." No, because the hydrogen atom has internal parts. It may, therefore, have various states of different internal energy, and describing the real nature requires more detail.

The question is: Does a proton have internal parts? Do we have to describe a proton by giving all possible states of protons, and mesons, and strange particles? We don't know. And even though we suppose that the electron is simple, so that all we have to tell about it is its momentum and its spin, maybe tomorrow we will discover that the electron also has inner gears and wheels. It would mean that our representation is incomplete, or wrong, or approximate—in the same way that a representation of the hydrogen atom which describes only its momentum would be incomplete, because it disregarded the fact that the hydrogen atom could have become excited inside. If an electron could become excited inside and turn into something else like, for instance, a muon, then it would be described not just by giving the states of the new particle, but presumably in terms of some more complicated internal wheels. The *main problem in the study of the fundamental particles today* is to discover what are the correct representations for the description of nature. At the present time, we *guess* that for the electron it is enough to specify its momentum and spin. We also guess that there is an idealized proton which has its π -mesons, and k -mesons, and so on, that all have to be specified. Several dozen particles—that's crazy! The question of what *is* a fundamental particle and what *is not* a fundamental particle—a subject you hear so much about these days—is the question of what is the final *representation* going to look like in the ultimate quantum mechanical description of the world. Will the electron's momentum still be the right thing with which to describe nature? Or even, should the whole question be put this way at all! This question must always come up in any scientific investigation. At any rate, we see a problem—how to find a representation. We don't know the answer. We don't even know whether we have the "right" problem, but if we do, we must first attempt to find out whether any particular particle is "fundamental" or not.

In the nonrelativistic quantum mechanics—if the energies are not too high, so that you don't disturb the inner workings of the strange particles and so forth—

you can do a pretty good job without worrying about these details. You can just decide to specify the momenta and spins of the electrons and of the nuclei; then everything will be all right. In most chemical reactions and other low-energy happenings, nothing goes on in the nuclei; they don't get excited. Furthermore, if a hydrogen atom is moving slowly and bumping quietly against other hydrogen atoms—never getting excited inside, or radiating, or anything complicated like that, but staying always in the ground state of energy for internal motion—you can use an approximation in which you talk about the hydrogen atom as one object, or particle, and not worry about the fact that it *can* do something inside. This will be a good approximation as long as the kinetic energy in any collision is well below 10 electron volts—the energy required to excite the hydrogen atom to a different internal state. We will often be making an approximation in which we do not include the possibility of inner motion, thereby decreasing the number of details that we have to put into our base states. Of course, we then omit some phenomena which would appear (usually) at some higher energy, but by making such approximations we can simplify very much the analysis of physical problems. For example, we can discuss the collision of two hydrogen atoms at low energy—or any chemical process—without worrying about the fact that the atomic nuclei could be excited. To summarize, then, when we can neglect the effects of any internal excited states of a particle we can choose a base set which are the states of definite momentum and z-component of angular momentum.

One problem then in describing nature is to find a suitable representation for the base states. But that's only the beginning. We still want to be able to say what "happens." If we know the "condition" of the world at one moment, we would like to know the condition at a later moment. So we also have to find the laws that determine how things change with time. We now address ourselves to this second part of the framework of quantum mechanics—how states change with time.

8-4 How states change with time

We have already talked about how we can represent a situation in which we put something through an apparatus. Now one convenient, delightful "apparatus" to consider is merely a wait of a few minutes; that is, you prepare a state ϕ , and then before you analyze it, you just let it sit. Perhaps you let it sit in some particular electric or magnetic field—it depends on the physical circumstances in the world. At any rate, whatever the conditions are, you let the object sit from time t_1 to time t_2 . Suppose that it is let out of your first apparatus in the condition ϕ at t_1 . And then it goes through an "apparatus," but the "apparatus" consists of just delay until t_2 . During the delay, various things could be going on—external forces applied or other shenanigans—so that something is happening. At the end of the delay, the amplitude to find the thing in some state χ is no longer exactly the same as it would have been without the delay. Since "waiting" is just a special case of an "apparatus," we can describe what happens by giving an amplitude with the same form as Eq. (8.17). Because the operation of "waiting" is especially important, we'll call it U instead of A , and to specify the starting and finishing times t_1 and t_2 , we'll write $U(t_2, t_1)$. The amplitude we want is

$$\langle \chi | U(t_2, t_1) | \phi \rangle. \quad (8.27)$$

Like any other such amplitude, it can be represented in some base system or other by writing it

$$\sum_{ij} \langle \chi | i \rangle \langle i | U(t_2, t_1) | j \rangle \langle j | \phi \rangle. \quad (8.28)$$

Then U is completely described by giving the whole set of amplitudes—the matrix

$$\langle i | U(t_2, t_1) | j \rangle. \quad (8.29)$$

We can point out, incidentally, that the matrix $\langle i | U(t_2, t_1) | j \rangle$ gives much more detail than may be needed. The high-class theoretical physicist working in

high-energy physics considers problems of the following general nature (because it's the way experiments are usually done). He starts with a couple of particles, like a proton and a proton, coming together from infinity. (In the lab, usually one particle is standing still, and the other comes from an accelerator that is practically at infinity on atomic level.) The things go crash and out come, say, two k -mesons, six π -mesons, and two neutrons in certain directions with certain momenta. What's the amplitude for this to happen? The mathematics looks like this: The ϕ -state specifies the spins and momenta of the incoming particles. The χ would be the question about what comes out. For instance, with what amplitude do you get the six mesons going in such-and-such directions, and the two neutrons going off in these directions, with their spins so-and-so. In other words, χ would be specified by giving all the momenta, and spins, and so on of the final products. Then the job of the theorist is to calculate the amplitude (8.27). However, he is really only interested in the special case that t_1 is $-\infty$ and t_2 is $+\infty$. (There is no experimental evidence on the details of the process, only on what comes in and what goes out.) The limiting case of $U(t_2, t_1)$ as $t_1 \rightarrow -\infty$ and $t_2 \rightarrow +\infty$ is called S , and what he wants is

$$\langle \chi | S | \phi \rangle.$$

Or, using the form (8.28), he would calculate the matrix

$$\langle i | S | j \rangle,$$

which is called the *S-matrix*. So if you see a theoretical physicist pacing the floor and saying, "All I have to do is calculate the *S-matrix*," you will know what he is worried about.

How to analyze—how to specify the laws for—the *S-matrix* is an interesting question. In relativistic quantum mechanics for high energies, it is done one way, but in nonrelativistic quantum mechanics it can be done another way, which is very convenient. (This other way can also be done in the relativistic case, but then it is not so convenient.) It is to work out the *U-matrix* for a small interval of time—in other words for t_2 and t_1 close together. If we can find a sequence of such *U*'s for successive intervals of time we can watch how things go as a function of time. You can appreciate immediately that this way is not so good for relativity, because you don't want to have to specify how everything looks "simultaneously" everywhere. But we won't worry about that—we're just going to worry about non-relativistic mechanics.

Suppose we think of the matrix *U* for a delay from t_1 until t_3 which is greater than t_2 . In other words, let's take three successive times: t_1 less than t_2 less than t_3 . Then we claim that the matrix that goes between t_1 and t_3 is the *product* in succession of what happens when you delay from t_1 until t_2 and then from t_2 until t_3 . It's just like the situation when we had two apparatuses *B* and *A* in series. We can then write, following the notation of Section 5-6,

$$U(t_3, t_1) = U(t_3, t_2) \cdot U(t_2, t_1). \quad (8.30)$$

In other words, we can analyze any time interval if we can analyze a sequence of short time intervals in between. We just multiply together all the pieces; that's the way that quantum mechanics is analyzed nonrelativistically.

Our problem, then, is to understand the matrix $U(t_2, t_1)$ for an infinitesimal time interval—for $t_2 = t_1 + \Delta t$. We ask ourselves this: If we have a state ϕ now, what does the state look like an infinitesimal time Δt later? Let's see how we write that out. Call the state at the time t , $|\psi(t)\rangle$ (we show the time dependence of ψ to be perfectly clear that we mean the condition at the time t). Now we ask the question: What is the condition after the small interval of time Δt later? The answer is

$$|\psi(t + \Delta t)\rangle = U(t + \Delta t, t) |\psi(t)\rangle. \quad (8.31)$$

This means the same as we meant by (8.25), namely, that the amplitude to

find x at the time $t + \Delta t$, is

$$\langle x | \psi(t + \Delta t) \rangle = \langle x | U(t + \Delta t, t) | \psi(t) \rangle. \quad (8.32)$$

Since we're not yet too good at these abstract things, let's project our amplitudes into a definite representation. If we multiply both sides of Eq. (8.31) by $\langle i |$, we get

$$\langle i | \psi(t + \Delta t) \rangle = \langle i | U(t + \Delta t, t) | \psi(t) \rangle. \quad (8.33)$$

We can also resolve the $|\psi(t)\rangle$ into base states and write

$$\langle i | \psi(t + \Delta t) \rangle = \sum_j \langle i | U(t + \Delta t, t) | j \rangle \langle j | \psi(t) \rangle. \quad (8.34)$$

We can understand Eq. (8.34) in the following way. If we let $C_i(t) = \langle i | \psi(t) \rangle$ stand for the amplitude to be in the base state i at the time t , then we can think of this amplitude (just a *number*, remember!) varying with time. Each C_i becomes a function of t . And we also have some information on *how* the amplitudes C_i vary with time. Each amplitude at $(t + \Delta t)$ is proportional to *all of the other* amplitudes at t multiplied by a set of coefficients. Let's call the U -matrix U_{ij} , by which we mean

$$U_{ij} = \langle i | U | j \rangle.$$

Then we can write Eq. (8.34) as

$$C_i(t + \Delta t) = \sum_j U_{ij}(t + \Delta t, t) C_j(t). \quad (8.35)$$

This, then, is how the dynamics of quantum mechanics is going to look.

We don't know much about the U_{ij} yet, except for one thing. We know that if Δt goes to zero, nothing can happen—we should get just the original state. So, $U_{ii} \rightarrow 1$ and $U_{ij} \rightarrow 0$, if $i \neq j$. In other words, $U_{ij} \rightarrow \delta_{ij}$ for $\Delta t \rightarrow 0$. Also, we can suppose that for small Δt , each of the coefficients U_{ij} should differ from δ_{ij} by amounts proportional to Δt ; so we can write

$$U_{ij} = \delta_{ij} + K_{ij} \Delta t. \quad (8.36)$$

However, it is usual to take the factor $(-i/\hbar)^\dagger$ out of the coefficients K_{ij} , for historical and other reasons; we prefer to write

$$U_{ij}(t + \Delta t, t) = \delta_{ij} - \frac{i}{\hbar} H_{ij}(t) \Delta t. \quad (8.37)$$

It is, of course, the same as Eq. (8.36) and, if you wish, just defines the coefficients $H_{ij}(t)$. The terms H_{ij} are just the derivatives with respect to t_2 of the coefficients $U_{ij}(t_2, t_1)$, evaluated at $t_2 = t_1 = t$.

Using this form for U in Eq. (8.35), we have

$$C_i(t + \Delta t) = \sum_j \left[\delta_{ij} - \frac{i}{\hbar} H_{ij}(t) \Delta t \right] C_j(t). \quad (8.38)$$

Taking the sum over the δ_{ij} term, we get just $C_i(t)$, which we can put on the other side of the equation. Then dividing by Δt , we have what we recognize as a derivative

$$\frac{C_i(t + \Delta t) - C_i(t)}{\Delta t} = -\frac{i}{\hbar} \sum_j H_{ij}(t) C_j(t)$$

or

$$i\hbar \frac{dC_i(t)}{dt} = \sum_j H_{ij}(t) C_j(t). \quad (8.39)$$

[†] We are in a bit of trouble here with notation. In the factor $(-i/\hbar)$, the i means the imaginary unit $\sqrt{-1}$, and *not* the index i that refers to the i th base state! We hope that you won't find it too confusing.

You remember that $C_i(t)$ is the amplitude $\langle i | \psi \rangle$ to find the state ψ in one of the base states i (at the time t). So Eq. (8.39) tells us how each of the coefficients $\langle i | \psi \rangle$ varies with time. But that is the same as saying that Eq. (8.39) tells us how the state ψ varies with time, since we are describing ψ in terms of the amplitudes $\langle i | \psi \rangle$. The variation of ψ in time is described in terms of the matrix H_{ij} , which has to include, of course, the things we are doing to the system to cause it to change. If we know the H_{ij} —which contains the physics of the situation and can, in general, depend on the time—we have a complete description of the behavior in time of the system. Equation (8.39) is then the quantum mechanical law for the dynamics of the world.

(We should say that we will always take a set of base states which are fixed and do not vary with time. There are people who use base states that also vary. However, that's like using a rotating coordinate system in mechanics, and we don't want to get involved in such complications.)

8-5 The Hamiltonian matrix

The idea, then, is that to describe the quantum mechanical world we need to pick a set of base states i and to write the physical laws by giving the matrix of coefficients H_{ij} . Then we have everything—we can answer any question about what will happen. So we have to learn what the rules are for finding the H 's to go with any physical situation—what corresponds to a magnetic field, or an electric field, and so on. And that's the hardest part. For instance, for the new strange particles, we have no idea what H_{ij} 's to use. In other words, no one knows the *complete* H_{ij} for the whole world. (Part of the difficulty is that one can hardly hope to discover the H_{ij} when no one even knows what the base states are!) We do have excellent approximations for nonrelativistic phenomena and for some other special cases. In particular, we have the forms that are needed for the motions of electrons in atoms—to describe chemistry. But we don't know the full true H for the whole universe.

The coefficients H_{ij} are called *the Hamiltonian matrix* or, for short, just *the Hamiltonian*. (How Hamilton, who worked in the 1830's, got his name on a quantum mechanical matrix is a tale of history.) It would be much better called the *energy matrix*, for reasons that will become apparent as we work with it. So *the* problem is: Know your Hamiltonian!

The Hamiltonian has one property that can be deduced right away, namely, that

$$H_{ij}^* = H_{ji}. \quad (8.40)$$

This follows from the condition that the total probability that the system is in *some* state does not change. If you start with a particle—an object or the world—then you've still got it as time goes on. The total probability of finding it somewhere is

$$\sum_i |C_i(t)|^2,$$

which must not vary with time. If this is to be true for any starting condition ϕ , then Eq. (8.40) must also be true.

As our first example, we take a situation in which the physical circumstances are not changing with time; we mean the *external* physical conditions, so that H is independent of time. Nobody is turning magnets on and off. We also pick a system for which only one base state is required for the description; it is an approximation we could make for a hydrogen atom at rest, or something similar. Equation (8.39) then says

$$i\hbar \frac{dC_1}{dt} = H_{11}C_1. \quad (8.41)$$

Only one equation—that's all! And if H_{11} is constant, this differential equation is easily solved to give

$$C_1 = (\text{const})e^{-(i/\hbar)H_{11}t} \quad (8.42)$$

This is the time dependence of a state with a definite energy $E = H_{11}$. You see why H_{ij} ought to be called the energy matrix. It is the generalization of the energy for more complex situations.

Next, to understand a little more about what the equations mean, we look at a system which has two base states. Then Eq. (8.39) reads

$$\begin{aligned} i\hbar \frac{dC_1}{dt} &= H_{11}C_1 + H_{12}C_2, \\ i\hbar \frac{dC_2}{dt} &= H_{21}C_1 + H_{22}C_2. \end{aligned} \quad (8.43)$$

If the H 's are again independent of time, you can easily solve these equations. We leave you to try for fun, and we'll come back and do them later. Yes, you can solve the quantum mechanics without knowing the H 's, so long as they are independent of time.

8-6 The ammonia molecule

We want now to show you how the dynamical equation of quantum mechanics can be used to describe a particular physical circumstance. We have picked an interesting but simple example in which, by making some reasonable guesses about the Hamiltonian, we can work out some important—and even practical—results. We are going to take a situation describable by two states: the ammonia molecule.

The ammonia molecule has one nitrogen atom and three hydrogen atoms located in a plane below the nitrogen so that the molecule has the form of a pyramid, as drawn in Fig. 8-1(a). Now this molecule, like any other, has an infinite number of states. It can spin around any possible axis; it can be moving in any direction; it can be vibrating inside, and so on, and so on. It is, therefore, not a two-state system at all. But we want to make an approximation that all other states remain fixed, because they don't enter into what we are concerned with at the moment. We will consider only that the molecule is spinning around its axis of symmetry (as shown in the figure), that it has zero translational momentum, and that it is vibrating as little as possible. That specifies all conditions except one: *there are still the two possible positions for the nitrogen atom*—the nitrogen may be on one side of the plane of hydrogen atoms or on the other, as shown in Fig. 8-1(a) and (b). So we will discuss the molecule as though it were a two-state system. We mean that there are only two states we are going to really worry about, all other things being assumed to stay put. You see, even if we know that it is spinning with a certain angular momentum around the axis and that it is moving with a certain momentum and vibrating in a definite way, there are still two possible states. We will say that the molecule is in the state $|1\rangle$ when the nitrogen is “up,” as in Fig. 8-1(a), and is in the state $|2\rangle$ when the nitrogen is “down,” as in (b). The states $|1\rangle$ and $|2\rangle$ will be taken as the set of base states for our analysis of the behavior of the ammonia molecule. At any moment, the actual state $|\psi\rangle$ of the molecule can be represented by giving $C_1 = \langle 1 | \psi \rangle$, the amplitude to be in state $|1\rangle$, and $C_2 = \langle 2 | \psi \rangle$, the amplitude to be in state $|2\rangle$. Then, using Eq. (8.8) we can write the state vector $|\psi\rangle$ as

$$|\psi\rangle = |1\rangle\langle 1 | \psi \rangle + |2\rangle\langle 2 | \psi \rangle$$

or

$$|\psi\rangle = |1\rangle C_1 + |2\rangle C_2. \quad (8.44)$$

Now the interesting thing is that if the molecule is known to be in some state at some instant, it will *not* be in the same state a little while later. The two C -coefficients will be changing with time according to the equations (8.43)—which hold for any two-state system. Suppose, for example, that you had made some observation—or had made some selection of the molecules—so that you *know* that the molecule is *initially* in the state $|1\rangle$. At some later time, there is some chance that it will be found in state $|2\rangle$. To find out what this chance is, we have to solve the differential equation which tells us how the amplitudes change with time.

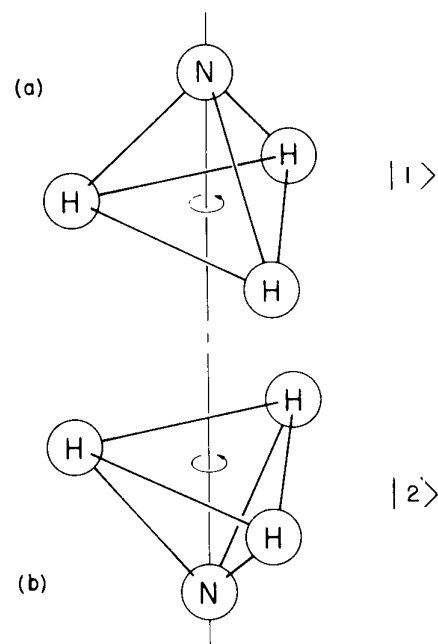


Fig. 8-1. Two equivalent geometric arrangements of the ammonia molecule.

The only trouble is that we don't know what to use for the coefficients H_{ij} in Eq. (8.43). There are some things we *can* say, however. Suppose that once the molecule was in the state $|1\rangle$ there was no chance that it could ever get into $|2\rangle$, and vice versa. Then H_{12} and H_{21} would both be zero, and Eq. (8.43) would read

$$i\hbar \frac{dC_1}{dt} = H_{11}C_1, \quad i\hbar \frac{dC_2}{dt} = H_{22}C_2.$$

We can easily solve these two equations; we get

$$C_1 = (\text{const})e^{-(i/\hbar)H_{11}t}, \quad C_2 = (\text{const})e^{-(i/\hbar)H_{22}t}. \quad (8.45)$$

These are just the amplitudes for *stationary* states with the energies $E_1 = H_{11}$ and $E_2 = H_{22}$. We note, however, that for the ammonia molecule the two states $|1\rangle$ and $|2\rangle$ have a definite symmetry. If nature is at all reasonable, the matrix elements H_{11} and H_{22} must be equal. We'll call them both E_0 , because they correspond to the energy the states would have if H_{12} and H_{21} were zero. But Eqs. (8.45) do not tell us what ammonia really does. It turns out that it is possible for the nitrogen to push its way through the three hydrogens and flip to the other side. It is quite difficult; to get half-way through requires a lot of energy. How can it get through if it hasn't got enough energy? There is *some* amplitude that it *will* penetrate the energy barrier. It is possible in quantum mechanics to sneak quickly across a region which is illegal energetically. There is, therefore, some small amplitude that a molecule which starts in $|1\rangle$ will get to the state $|2\rangle$. The coefficients H_{12} and H_{21} are not really zero. Again, by symmetry, they should both be the same—at least in magnitude. In fact, we already know that, in general, H_{ij} must be equal to the complex conjugate of H_{ji} , so they can differ only by a phase. It turns out, as you will see, that there is no loss of generality if we take them equal to each other. For later convenience we set them equal to a negative number; we take $H_{12} = H_{21} = -A$. We then have the following pair of equations:

$$i\hbar \frac{dC_1}{dt} = E_0C_1 - AC_2, \quad (8.46)$$

$$i\hbar \frac{dC_2}{dt} = E_0C_2 - AC_1. \quad (8.47)$$

These equations are simple enough and can be solved in any number of ways. One convenient way is the following. Taking the sum of the two, we get

$$i\hbar \frac{d}{dt} (C_1 + C_2) = (E_0 - A)(C_1 + C_2),$$

whose solution is

$$C_1 + C_2 = ae^{-(i/\hbar)(E_0-A)t}. \quad (8.48)$$

Then, taking the difference of (8.46) and (8.47), we find that

$$i\hbar \frac{d}{dt} (C_1 - C_2) = (E_0 + A)(C_1 - C_2),$$

which gives

$$C_1 - C_2 = be^{-(i/\hbar)(E_0+A)t}. \quad (8.49)$$

We have called the two integration constants a and b ; they are, of course, to be chosen to give the appropriate starting condition for any particular physical problem. Now, by adding and subtracting (8.48) and (8.49), we get C_1 and C_2 :

$$C_1(t) = \frac{a}{2} e^{-(i/\hbar)(E_0-A)t} + \frac{b}{2} e^{-(i/\hbar)(E_0+A)t}, \quad (8.50)$$

$$C_2(t) = \frac{a}{2} e^{-(i/\hbar)(E_0-A)t} - \frac{b}{2} e^{-(i/\hbar)(E_0+A)t}. \quad (8.51)$$

They are the same except for the sign of the second term.

We have the solutions; now what do they mean? (The trouble with quantum mechanics is not only in solving the equations but in understanding what the solutions mean!) First, notice that if $b = 0$, both terms have the same frequency $\omega = (E_0 - A)/\hbar$. If everything changes at one frequency, it means that the system is in a state of definite energy—here, the energy $(E_0 - A)$. So there is a stationary state of this energy in which the two amplitudes C_1 and C_2 are equal. We get the result that *the ammonia molecule has a definite energy* $(E_0 - A)$ if there are equal amplitudes for the nitrogen atom to be “up” and to be “down.”

There is another stationary state possible if $a = 0$; both amplitudes then have the frequency $(E_0 + A)/\hbar$. So there is another state with the definite energy $(E_0 + A)$ if the two amplitudes are equal but with the opposite sign; $C_2 = -C_1$. These are the only two states of definite energy. We will discuss the states of the ammonia molecule in more detail in the next chapter; we will mention here only a couple of things.

We conclude that *because* there is some chance that the nitrogen atom can flip from one position to the other, the energy of the molecule is not just E_0 , as we would have expected, but that there are *two* energy levels $(E_0 + A)$ and $(E_0 - A)$. Every one of the possible states of the molecule, whatever energy it has, is “split” into two levels. We say *every* one of the states because, you remember, we picked out one particular state of rotation, and internal energy, and so on. For each possible condition of that kind there is a doublet of energy levels because of the flip-flop of the molecule.

Let’s now ask the following question about an ammonia molecule. Suppose that at $t = 0$, we *know* that a molecule is in the state $|1\rangle$ or, in other words, that $C_1(0) = 1$ and $C_2(0) = 0$. What is the probability that the molecule will be found in the state $|2\rangle$ at the time t , or will still be found in state $|1\rangle$ at the time t ? Our starting condition tells us what a and b are in Eqs. (8.50) and (8.51). Letting $t = 0$, we have that

$$C_1(0) = \frac{a+b}{2} = 1, \quad C_2(0) = \frac{a-b}{2} = 0.$$

Clearly, $a = b = 1$. Putting these values into the formulas for $C_1(t)$ and $C_2(t)$ and rearranging some terms, we have

$$C_1(t) = e^{-(i/\hbar)E_0 t} \left(\frac{e^{(i/\hbar)At} + e^{-(i/\hbar)At}}{2} \right),$$

$$C_2(t) = e^{-(i/\hbar)E_0 t} \left(\frac{e^{(i/\hbar)At} - e^{-(i/\hbar)At}}{2} \right).$$

We can rewrite these as

$$C_1(t) = e^{-(i/\hbar)E_0 t} \cos \frac{At}{\hbar}, \quad (8.52)$$

$$C_2(t) = ie^{-(i/\hbar)E_0 t} \sin \frac{At}{\hbar}. \quad (8.53)$$

The two amplitudes have a magnitude that varies harmonically with time.

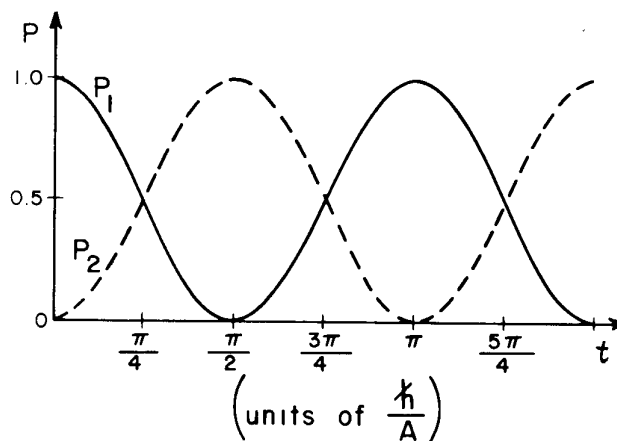
The probability that the molecule is found in state $|2\rangle$ at the time t is the absolute square of $C_2(t)$:

$$|C_2(t)|^2 = \sin^2 \frac{At}{\hbar}. \quad (8.54)$$

The probability starts at zero (as it should), rises to one, and then oscillates back and forth between zero and one, as shown in the curve marked P_2 of Fig. 8-2. The probability of being in the $|1\rangle$ state does not, of course, stay at one. It “dumps” into the second state until the probability of finding the molecule in the first state is zero, as shown by the curve P_1 of Fig. 8-2. The probability sloshes back and forth between the two.

A long time ago we saw what happens when we have two equal pendulums with a slight coupling. (See Chapter 49, Vol. I.) When we lift one back and let go,

Fig. 8-2. The probability P_1 that an ammonia molecule in state $|1\rangle$ at $t = 0$ will be found in state $|1\rangle$ at t . The probability P_2 that it will be found in state $|2\rangle$.



it swings, but then gradually the other one starts to swing. Pretty soon the second pendulum has picked up all the energy. Then, the process reverses, and pendulum number one picks up the energy. It is exactly the same kind of a thing. The speed at which the energy is swapped back and forth depends on the coupling between the two pendulums—the rate at which the “oscillation” is able to leak across. Also, you remember, with the two pendulums there are two special motions—each with a definite frequency—which we call the fundamental modes. If we pull both pendulums out together, they swing together at one frequency. On the other hand, if we pull one out one way and the other out the other way, there is another stationary mode also at a definite frequency.

Well, here we have a similar situation—the ammonia molecule is mathematically like the pair of pendulums. These are the two frequencies— $(E_0 + A)/\hbar$ and $(E_0 - A)/\hbar$ —for when they are oscillating together, or oscillating opposite.

The pendulum analogy is not much deeper than the principle that the same equations have the same solutions. The linear equations for the amplitudes (8.39) are very much like the linear equations of harmonic oscillators. (In fact, this is the reason behind the success of our classical theory of the index of refraction, in which we replaced the quantum mechanical atom by a harmonic oscillator, even though, classically, this is not a reasonable view of electrons circulating about a nucleus.) If you pull the nitrogen to one side, then you get a *superposition* of these two frequencies, and you get a kind of *beat* note, because the system is *not* in one or the other states of definite frequency. The splitting of the energy levels of the ammonia molecule is, however, strictly a quantum mechanical effect.

The splitting of the energy levels of the ammonia molecule has important practical applications which we will describe in the next chapter. At long last we have an example of a practical physical problem that you can understand with the quantum mechanics!

The Ammonia Maser

9-1 The states of an ammonia molecule

In this chapter we are going to discuss the application of quantum mechanics to a practical device, the ammonia maser. You may wonder why we stop our formal development of quantum mechanics to do a special problem, but you will find that many of the features of this special problem are quite common in the general theory of quantum mechanics, and you will learn a great deal by considering this one problem in detail. The ammonia maser is a device for generating electromagnetic waves, whose operation is based on the properties of the ammonia molecule which we discussed briefly in the last chapter. We begin by summarizing what we found there.

The ammonia molecule has many states, but we are considering it as a two-state system, thinking now only about what happens when the molecule is in any specific state of rotation or translation. A physical model for the two states can be visualized as follows. If the ammonia molecule is considered to be rotating about an axis passing through the nitrogen atom and perpendicular to the plane of the hydrogen atoms, as shown in Fig. 9-1, there are still two possible conditions—the nitrogen may be on one side of the plane of hydrogen atoms or on the other. We call these two states $|1\rangle$ and $|2\rangle$. They are taken as a set of base states for our analysis of the behavior of the ammonia molecule.

9-1 The states of an ammonia molecule

9-2 The molecule in a static electric field

9-3 Transitions in a time-dependent field

9-4 Transitions at resonance

9-5 Transitions off resonance

9-6 The absorption of light

MASER = Microwave Amplification by Stimulated Emission of Radiation

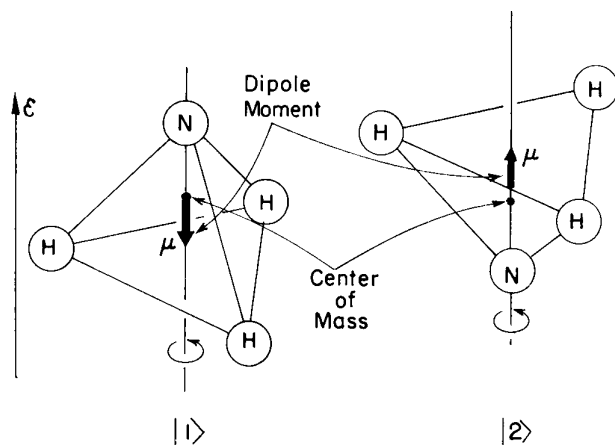


Fig. 9-1. A physical model of two base states for the ammonia molecule. These states have the electric dipole moments μ .

In a system with two base states, any state $|\psi\rangle$ of the system can always be described as a linear combination of the two base states; that is, there is a certain amplitude C_1 to be in one base state and an amplitude C_2 to be in the other. We can write its state vector as

$$|\psi\rangle = |1\rangle C_1 + |2\rangle C_2, \quad (9.1)$$

where

$$C_1 = \langle 1 | \psi \rangle \quad \text{and} \quad C_2 = \langle 2 | \psi \rangle.$$

These two amplitudes change with time according to the Hamiltonian equations, Eq. (8.43). Making use of the symmetry of the two states of the ammonia molecule, we set $H_{11} = H_{22} = E_0$, and $H_{12} = H_{21} = -A$, and get the

solution [see Eqs. (8.50) and (8.51)]

$$C_1 = \frac{a}{2} e^{-(i/\hbar)(E_0-A)t} + \frac{b}{2} e^{-(i/\hbar)(E_0+A)t}, \quad (9.2)$$

$$C_2 = \frac{a}{2} e^{-(i/\hbar)(E_0-A)t} - \frac{b}{2} e^{-(i/\hbar)(E_0+A)t}. \quad (9.3)$$

We want now to take a closer look at these general solutions. Suppose that the molecule was initially put into a state $|\psi_{II}\rangle$ for which the coefficient b was equal to zero. Then at $t = 0$ the amplitudes to be in the states $|I\rangle$ and $|2\rangle$ are identical, *and they stay that way for all time*. Their phases both vary with time in the same way—with the frequency $(E_0 - A)/\hbar$. Similarly, if we were to put the molecule into a state $|\psi_I\rangle$ for which $a = 0$, the amplitude C_2 is the negative of C_1 , and this relationship would stay that way forever. Both amplitudes would now vary with time with the frequency $(E_0 + A)/\hbar$. These are the only two possibilities of states for which the relation between C_1 and C_2 is independent of time.

We have found two special solutions in which the two amplitudes *do not vary in magnitude* and, furthermore, have phases which vary at the same frequencies. These are *stationary states* as we defined them in Section 7-1, which means that they are *states of definite energy*. The state $|\psi_{II}\rangle$ has the energy $E_{II} = E_0 - A$, and the state $|\psi_I\rangle$ has the energy $E_I = E_0 + A$. They are the only two stationary states that exist, so we find that the molecule has two energy levels, with the energy difference $2A$. (We mean, of course, two energy levels for the assumed state of rotation and vibration which we referred to in our initial assumptions.)†

If we hadn't allowed for the possibility of the nitrogen flipping back and forth, we would have taken A equal to zero and the two energy levels would be on top of each other at energy E_0 . The actual levels are not this way; their *average* energy is E_0 , but they are split apart by $\pm A$, giving a separation of $2A$ between the energies of the two states. Since A is, in fact, very small, the difference in energy is also very small.

In order to excite an *electron* inside an atom, the energies involved are relatively very high—requiring photons in the optical or ultraviolet range. To excite the *vibrations* of the molecules involves photons in the infrared. If you talk about exciting *rotations*, the energy differences of the states correspond to photons in the far infrared. But the energy difference $2A$ is lower than any of those and is, in fact, below the infrared and well into the microwave region. Experimentally, it has been found that there is a pair of energy levels with a separation of 10^{-4} electron volt—corresponding to a frequency 24,000 megacycles. Evidently this means that $2A = hf$, with $f = 24,000$ megacycles (corresponding to a wavelength of $1\frac{1}{4}$ cm). So here we have a molecule that has a transition which does not emit light in the ordinary sense, but emits microwaves.

For the work that follows we need to describe these two states of definite energy a little bit better. Suppose we were to construct an amplitude C_{II} by taking the sum of the two numbers C_1 and C_2 :

$$C_{II} = C_1 + C_2 = \langle I | \Phi \rangle + \langle 2 | \Phi \rangle. \quad (9.4)$$

What would that mean? Well, this is just the amplitude to find the state $|\Phi\rangle$ in a new state $|II\rangle$ in which the amplitudes of the original base states are equal. That is, writing $C_{II} = \langle II | \Phi \rangle$, we can abstract the $|\Phi\rangle$ away from Eq. (9.4)—because it is true for any Φ —and get

$$\langle II | = \langle I | + \langle 2 |,$$

which means the same as

$$|II\rangle = |I\rangle + |2\rangle. \quad (9.5)$$

† In what follows it is helpful—in reading to yourself or in talking to someone else—to have a handy way of distinguishing between the Arabic 1 and 2 and the Roman I and II. We find it convenient to reserve the names “one” and “two” for the Arabic numbers, and to call I and II by the names “eins” and “zwei” (although “unus” and “duo” might be more logical!).

The amplitude for the state $|II\rangle$ to be in the state $|I\rangle$ is

$$\langle I | II \rangle = \langle I | I \rangle + \langle I | 2 \rangle,$$

which is, of course, just 1, since $|I\rangle$ and $|2\rangle$ are base states. The amplitude for the state $|II\rangle$ to be in the state $|2\rangle$ is also 1, so the state $|II\rangle$ is one which has equal amplitudes to be in the two base states $|I\rangle$ and $|2\rangle$.

We are, however, in a bit of trouble. The state $|II\rangle$ has a total probability greater than one of being in *some* base state *or other*. That simply means, however, that the state vector is not properly “normalized.” We can take care of that by remembering that we should have $\langle II | II \rangle = 1$, which must be so for any state. Using the general relation that

$$\langle \chi | \Phi \rangle = \sum_i \langle \chi | i \rangle \langle i | \Phi \rangle,$$

letting both Φ and χ be the state II , and taking the sum over the base states $|I\rangle$ and $|2\rangle$, we get that

$$\langle II | II \rangle = \langle II | I \rangle \langle I | II \rangle + \langle II | 2 \rangle \langle 2 | II \rangle.$$

This will be equal to one as it should if we change our definition of C_{II} —in Eq. (9.4)—to read

$$C_{II} = \frac{1}{\sqrt{2}} [C_1 + C_2].$$

In the same way we can construct an amplitude

$$C_I = \frac{1}{\sqrt{2}} [C_1 - C_2],$$

or

$$C_I = \frac{1}{\sqrt{2}} [\langle I | \Phi \rangle - \langle 2 | \Phi \rangle]. \quad (9.6)$$

This amplitude is the projection of the state $|\Phi\rangle$ into a new state $|I\rangle$ which has opposite amplitudes to be in the states $|I\rangle$ and $|2\rangle$. Namely, Eq. (9.6) means the same as

$$\langle I | = \frac{1}{\sqrt{2}} [\langle I | - \langle 2 |],$$

or

$$|I\rangle = \frac{1}{\sqrt{2}} [|I\rangle - |2\rangle], \quad (9.7)$$

from which it follows that

$$\langle I | I \rangle = \frac{1}{\sqrt{2}} = -\langle 2 | I \rangle.$$

Now the reason we have done all this is that the states $|I\rangle$ and $|II\rangle$ *can be taken as a new set of base states* which are especially convenient for describing the stationary states of the ammonia molecule. You remember that the requirement for a set of base states is that

$$\langle i | j \rangle = \delta_{ij}.$$

We have already fixed things so that

$$\langle I | I \rangle = \langle II | II \rangle = 1.$$

You can easily show from Eqs. (9.5) and (9.7) that

$$\langle I | II \rangle = \langle II | I \rangle = 0.$$

The amplitudes $C_I = \langle I | \Phi \rangle$ and $C_{II} = \langle II | \Phi \rangle$ for any state Φ to be in our new base states $|I\rangle$ and $|II\rangle$ must also satisfy a Hamiltonian equation with the

form of Eq. (8.39). In fact, if we just subtract the two equations (9.2) and (9.3) and differentiate with respect to t , we see that

$$i\hbar \frac{dC_I}{dt} = (E_0 + A)C_I = E_I C_I. \quad (9.8)$$

And taking the sum of Eqs. (9.2) and (9.3), we see that

$$i\hbar \frac{dC_{II}}{dt} = (E_0 - A)C_{II} = E_{II} C_{II}. \quad (9.9)$$

Using $|I\rangle$ and $|II\rangle$ for base states, the Hamiltonian matrix has the simple form

$$\begin{aligned} H_{I,I} &= E_I, & H_{I,II} &= 0, \\ H_{II,I} &= 0, & H_{II,II} &= E_{II}. \end{aligned}$$

Note that each of the Eqs. (9.8) and (9.9) look just like what we had in Section 8-6 for the equation of a one-state system. They have a simple exponential time dependence corresponding to a single energy. As time goes on, the amplitudes to be in each state act independently.

The two stationary states $|\psi_I\rangle$ and $|\psi_{II}\rangle$ we found above are, of course, solutions of Eqs. (9.8) and (9.9). The state $|\psi_I\rangle$ (for which $C_1 = -C_2$) has

$$C_I = e^{-(i/\hbar)(E_0+A)t}, \quad C_{II} = 0. \quad (9.10)$$

And the state $|\psi_{II}\rangle$ (for which $C_1 = C_2$) has

$$C_I = 0, \quad C_{II} = e^{-(i/\hbar)(E_0-A)t}. \quad (9.11)$$

Remember that the amplitudes in Eq. (9.10) are

$$C_I = \langle I | \psi_I \rangle, \quad \text{and} \quad C_{II} = \langle II | \psi_I \rangle;$$

so Eq. (9.10) means the same thing as

$$|\psi_I\rangle = |I\rangle e^{-(i/\hbar)(E_0+A)t}.$$

That is, the state vector of the stationary state $|\psi_I\rangle$ is the same as the state vector of the base state $|I\rangle$ except for the exponential factor appropriate to the energy of the state. In fact at $t = 0$

$$|\psi_I\rangle = |I\rangle;$$

the state $|I\rangle$ has the same physical configuration as the stationary state of energy $E_0 + A$. In the same way, we have for the second stationary state that

$$|\psi_{II}\rangle = |II\rangle e^{-(i/\hbar)(E_0-A)t}.$$

The state $|II\rangle$ is just the stationary state of energy $E_0 - A$ at $t = 0$. Thus our two new base states $|I\rangle$ and $|II\rangle$ have physically the form of the states of definite energy, with the exponential time factor taken out so that they can be time-independent base states. (In what follows we will find it convenient not to have to distinguish always between the stationary states $|\psi_I\rangle$ and $|\psi_{II}\rangle$ and their base states $|I\rangle$ and $|II\rangle$, since they differ only by the obvious time factors.)

In summary, the state vectors $|\psi_I\rangle$ and $|\psi_{II}\rangle$ are a pair of base vectors which are appropriate for describing the definite energy states of the ammonia molecule. They are related to our original base vectors by

$$|I\rangle = \frac{1}{\sqrt{2}} [|1\rangle - |2\rangle], \quad |II\rangle = \frac{1}{\sqrt{2}} [|1\rangle + |2\rangle]. \quad (9.12)$$

The amplitudes to be in $|I\rangle$ and $|II\rangle$ are related to C_1 and C_2 by

$$C_I = \frac{1}{\sqrt{2}} [C_1 - C_2], \quad C_{II} = \frac{1}{\sqrt{2}} [C_1 + C_2]. \quad (9.13)$$

Any state at all can be represented by a linear combination of $|I\rangle$ and $|2\rangle$ —with the coefficients C_1 and C_2 —or by a linear combination of the definite energy base states $|I\rangle$ and $|II\rangle$ —with the coefficients C_I and C_{II} . Thus,

$$|\Phi\rangle = |I\rangle C_1 + |2\rangle C_2$$

or

$$|\Phi\rangle = |I\rangle C_I + |II\rangle C_{II}.$$

The second form gives us the amplitudes for finding the state $|\Phi\rangle$ in a state with the energy $E_I = E_0 + A$ or in a state with the energy $E_{II} = E_0 - A$.

9-2 The molecule in a static electric field

If the ammonia molecule is in either of the two states of definite energy and we disturb it at a frequency ω such that $\hbar\omega = E_I - E_{II} = 2A$, the system may make a transition from one state to the other. Or, if it is in the upper state, it may change to the lower state and emit a photon. But in order to induce such transitions you must have a physical connection to the states—some way of disturbing the system. There must be some external machinery for affecting the states, such as magnetic or electric fields. In this particular case, these states are sensitive to an electric field. We will, therefore, look next at the problem of the behavior of the ammonia molecule in an external electric field.

To discuss the behavior in an electric field, we will go back to the original base system $|I\rangle$ and $|2\rangle$, rather than using $|I\rangle$ and $|II\rangle$. Suppose that there is an electric field in a direction perpendicular to the plane of the hydrogen atoms. Disregarding for the moment the possibility of flipping back and forth, would it be true that the energy of this molecule is the same for the two positions of the nitrogen atom? Generally, no. The electrons tend to lie closer to the nitrogen than to the hydrogen nuclei, so the hydrogens are slightly positive. The actual amount depends on the details of electron distribution. It is a complicated problem to figure out exactly what this distribution is, but in any case the net result is that the ammonia molecule has an electric dipole moment, as indicated in Fig. 9-1. We can continue our analysis without knowing in detail the direction or amount of displacement of the charge. However, to be consistent with the notation of others, let's suppose that the electric dipole moment is μ , with its direction point *from* the nitrogen atom and perpendicular to the plane of the hydrogen atoms.

Now, when the nitrogen flips from one side to the other, the center of mass will not move, but the electric dipole moment will flip over. As a result of this moment, the energy in an electric field \mathcal{E} will depend on the molecular orientation.† With the assumption made above, the potential energy will be higher if the nitrogen atom points in the direction of the field, and lower if it is in the opposite direction; the separation in the two energies will be $2\mu\mathcal{E}$.

In the discussion up to this point, we have assumed values of E_0 and A without knowing how to calculate them. According to the correct physical theory, it should be possible to calculate these constants in terms of the positions and motions of all the nuclei and electrons. But nobody has ever done it. Such a system involves ten electrons and four nuclei and that's just too complicated a problem. As a matter of fact, there is no one who knows much more about this molecule than we do. All anyone can say is that when there is an electric field, the energy of the two states is different, the difference being proportional to the electric field. We have called the coefficient of proportionality 2μ , but its value must be determined experimentally. We can also say that the molecule has the amplitude A to flip over, but this will have to be measured experimentally. Nobody can give us accurate theoretical values of μ and A , because the calculations are too complicated to do in detail.

† We are sorry that we have to introduce a new notation. Since we have been using p and E for momentum and energy, we don't want to use them again for dipole moment and electric field. Remember, in this section μ is the *electric* dipole moment.

For the ammonia molecule in an electric field, our description must be changed. If we ignored the amplitude for the molecule to flip from one configuration to the other, we would expect the energies of the two states $|1\rangle$ and $|2\rangle$ to be $(E_0 \pm \mu\mathcal{E})$. Following the procedure of the last chapter, we take

$$H_{11} = E_0 + \mu\mathcal{E}, \quad H_{22} = E_0 - \mu\mathcal{E}. \quad (9.14)$$

Also we will assume that for the electric fields of interest the field does not affect appreciably the geometry of the molecule and, therefore, does not affect the amplitude that the nitrogen will jump from one position to the other. We can then take that H_{12} and H_{21} are not changed; so

$$H_{12} = H_{21} = -A. \quad (9.15)$$

We must now solve the Hamiltonian equations, Eq. (8.43), with these new values of H_{ij} . We could solve them just as we did before, but since we are going to have several occasions to want the solutions for two-state systems, let's solve the equations once and for all in the general case of arbitrary H_{ij} —assuming only that they do not change with time.

We want the general solution of the pair of Hamiltonian equations

$$i\hbar \frac{dC_1}{dt} = H_{11}C_1 + H_{12}C_2, \quad (9.16)$$

$$i\hbar \frac{dC_2}{dt} = H_{21}C_1 + H_{22}C_2. \quad (9.17)$$

Since these are linear differential equations with constant coefficients, we can always find solutions which are exponential functions of the dependent variable t . We will first look for a solution in which C_1 and C_2 both have the same time dependence; we can use the trial functions

$$C_1 = a_1 e^{-i\omega t}, \quad C_2 = a_2 e^{-i\omega t}.$$

Since such a solution corresponds to a state of energy $E = \hbar\omega$, we may as well write right away

$$C_1 = a_1 e^{-(i/\hbar)Et}, \quad (9.18)$$

$$C_2 = a_2 e^{-(i/\hbar)Et}, \quad (9.19)$$

where E is as yet unknown and to be determined so that the differential equations (9.16) and (9.17) are satisfied.

When we substitute C_1 and C_2 from (9.18) and (9.19) in the differential equations (9.16) and (9.17), the derivatives give us just $-iE/\hbar$ times C_1 or C_2 , so the left sides become just EC_1 and EC_2 . Cancelling the common exponential factors, we get

$$Ea_1 = H_{11}a_1 + H_{12}a_2, \quad Ea_2 = H_{21}a_1 + H_{22}a_2.$$

Or, rearranging the terms, we have

$$(E - H_{11})a_1 - H_{12}a_2 = 0, \quad (9.20)$$

$$-H_{21}a_1 + (E - H_{22})a_2 = 0. \quad (9.21)$$

With such a set of homogeneous algebraic equations, there will be nonzero solutions for a_1 and a_2 only if the determinant of the coefficients of a_1 and a_2 is zero, that is, if

$$\text{Det} \begin{pmatrix} E - H_{11} & -H_{12} \\ -H_{21} & E - H_{22} \end{pmatrix} = 0. \quad (9.22)$$

However, when there are only two equations and two unknowns, we don't need such a sophisticated idea. The two equations (9.20) and (9.21) each give a ratio for the two coefficients a_1 and a_2 , and these two ratios must be equal. From (9.20) we have that

$$\frac{a_1}{a_2} = \frac{H_{12}}{E - H_{11}}, \quad (9.23)$$

and from (9.21) that

$$\frac{a_1}{a_2} = \frac{E - H_{22}}{H_{21}}. \quad (9.24)$$

Equating these two ratios, we get that E must satisfy

$$(E - H_{11})(E - H_{22}) - H_{12}H_{21} = 0.$$

This is the same result we would get by solving Eq. (9.22). Either way, we have a quadratic equation for E which has two solutions:

$$E = \frac{H_{11} + H_{22}}{2} \pm \sqrt{\frac{(H_{11} - H_{22})^2}{4} + H_{12}H_{21}}. \quad (9.25)$$

There are two possible values for the energy E . Note that both solutions give *real numbers* for the energy, because H_{11} and H_{22} are real, and $H_{12}H_{21}$ is equal to $H_{12}H_{12}^* = |H_{12}|^2$, which is both real and positive.

Using the same convention we took before, we will call the upper energy E_I and the lower energy E_{II} . We have

$$E_I = \frac{H_{11} + H_{22}}{2} + \sqrt{\frac{(H_{11} - H_{22})^2}{4} + H_{12}H_{21}}, \quad (9.26)$$

$$E_{II} = \frac{H_{11} + H_{22}}{2} - \sqrt{\frac{(H_{11} - H_{22})^2}{4} + H_{12}H_{21}}. \quad (9.27)$$

Using each of these two energies separately in Eqs. (9.18) and (9.19), we have the amplitudes for the two stationary states (the states of definite energy). If there are no external disturbances, a system initially in one of these states will stay that way forever—only its phase changes.

We can check our results for two special cases. If $H_{12} = H_{21} = 0$, we have that $E_I = H_{11}$ and $E_{II} = H_{22}$. This is certainly correct, because then Eqs. (9.16) and (9.17) are uncoupled, and each represents a state of energy H_{11} and H_{22} . Next, if we set $H_{11} = H_{22} = E_0$ and $H_{21} = H_{12} = -A$, we get the solution we found before:

$$E_I = E_0 + A \quad \text{and} \quad E_{II} = E_0 - A.$$

For the general case, the two solutions E_I and E_{II} refer to two states—which we can again call the states

$$|\psi_I\rangle = |I\rangle e^{-(i/\hbar)E_I t} \quad \text{and} \quad |\psi_{II}\rangle = |II\rangle e^{-(i/\hbar)E_{II} t}.$$

These states will have C_1 and C_2 as given in Eqs. (9.18) and (9.19), where a_1 and a_2 are still to be determined. Their ratio is given by either Eq. (9.23) or Eq. (9.24). They must also satisfy one more condition. If the system is known to be in one of the stationary states, the sum of the probabilities that it will be found in $|I\rangle$ or $|2\rangle$ must equal one. We must have that

$$|C_1|^2 + |C_2|^2 = 1, \quad (9.28)$$

or, equivalently,

$$|a_1|^2 + |a_2|^2 = 1. \quad (9.29)$$

These conditions do not uniquely specify a_1 and a_2 ; they are still undetermined

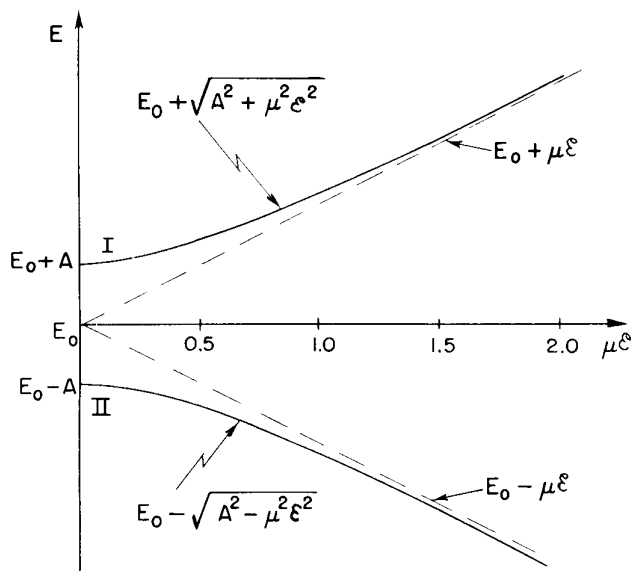


Fig. 9-2. Energy levels of the ammonia molecule in an electric field.

by an arbitrary phase—in other words, by a factor like $e^{i\delta}$. Although general solutions for the a 's can be written down,† it is usually more convenient to work them out for each special case.

Let's go back now to our particular example of the ammonia molecule in an electric field. Using the values for H_{11} , H_{22} , and H_{12} given in (9.14) and (9.15), we get for the energies of the two stationary states

$$E_I = E_0 + \sqrt{A^2 + \mu^2 \mathcal{E}^2}, \quad E_{II} = E_0 - \sqrt{A^2 + \mu^2 \mathcal{E}^2}. \quad (9.30)$$

These two energies are plotted as a function of the electric field strength \mathcal{E} in Fig. 9-2. When the electric field is zero, the two energies are, of course, just $E_0 \pm A$. When an electric field is applied, the splitting between the two levels increases. The splitting increases at first slowly with \mathcal{E} , but eventually becomes proportional to \mathcal{E} . (The curve is a hyperbola.) For enormously strong fields, the energies are just

$$E_I = E_0 + \mu \mathcal{E} = H_{11}, \quad E_{II} = E_0 - \mu \mathcal{E} = H_{22}. \quad (9.31)$$

The fact that there is an amplitude for the nitrogen to flip back and forth has little effect when the two positions have very different energies. This is an interesting point which we will come back to again later.

We are at last ready to understand the operation of the ammonia maser. The idea is the following. First, we find a way of separating molecules in the state $|I\rangle$ from those in the state $|II\rangle$.‡ Then the molecules in the higher energy state $|I\rangle$ are passed through a cavity which has a resonant frequency of 24,000 megacycles. The molecules can deliver energy to the cavity—in a way we will discuss later—and leave the cavity in the state $|II\rangle$. Each molecule that makes such a transition will deliver the energy $E = E_I - E_{II}$ to the cavity. The energy from the molecules will appear as electrical energy in the cavity.

How can we separate the two molecular states? One method is as follows. The ammonia gas is let out of a little jet and passed through a pair of slits to give a narrow beam, as shown in Fig. 9-3. The beam is then set through a

† For example, the following set is one acceptable solution, as you can easily verify:

$$a_1 = \frac{H_{12}}{[(E - H_{11})^2 + H_{12}H_{21}]^{1/2}}, \quad a_2 = \frac{E - H_{11}}{[(E - H_{11})^2 + H_{12}H_{21}]^{1/2}}.$$

‡ From now on we will write $|I\rangle$ and $|II\rangle$ instead of $|\psi_I\rangle$ and $|\psi_{II}\rangle$. You must remember that the actual states $|\psi_I\rangle$ and $|\psi_{II}\rangle$ are the energy base states multiplied by the appropriate exponential factor.

region in which there is a large transverse electric field. The electrodes to produce the field are shaped so that the electric field varies rapidly across the beam. Then the square of the electric field $\mathcal{E} \cdot \mathcal{E}$ will have a large gradient perpendicular to the beam. Now a molecule in state $|I\rangle$ has an energy which increases with \mathcal{E}^2 , and therefore this part of the beam will be deflected toward the region of lower \mathcal{E}^2 . A molecule in state $|II\rangle$ will, on the other hand, be deflected toward the region of larger \mathcal{E}^2 , since its energy decreases as \mathcal{E}^2 increases.

Incidentally, with the electric fields which can be generated in the laboratory, the energy $\mu\mathcal{E}$ is always much smaller than A . In such cases, the square root in Eqs. (9.30) can be approximated by

$$A \left(1 + \frac{1}{2} \frac{\mu^2 \mathcal{E}^2}{A^2} \right). \quad (9.32)$$

So the energy levels are, for all practical purposes,

$$E_I = E_0 + A + \frac{\mu^2 \mathcal{E}^2}{2A} \quad (9.33)$$

and

$$E_{II} = E_0 - A - \frac{\mu^2 \mathcal{E}^2}{2A}. \quad (9.34)$$

And the energies vary approximately linearly with \mathcal{E}^2 . The force on the molecules is then

$$\mathbf{F} = \frac{\mu^2}{2A} \nabla \mathcal{E}^2. \quad (9.35)$$

Many molecules have an energy in an electric field which is proportional to \mathcal{E}^2 . The coefficient is the polarizability of the molecule. Ammonia has an unusually high polarizability because of the small value of A in the denominator. Thus, ammonia molecules are unusually sensitive to an electric field. (What would you expect for the dielectric coefficient of NH_3 gas?)

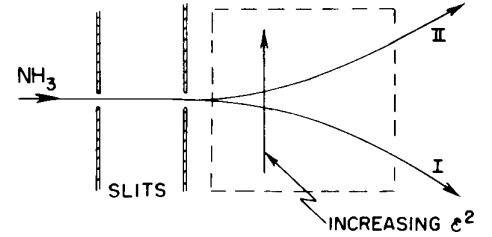


Fig. 9-3. The ammonia beam may be separated by an electric field in which \mathcal{E}^2 has a gradient perpendicular to the beam.

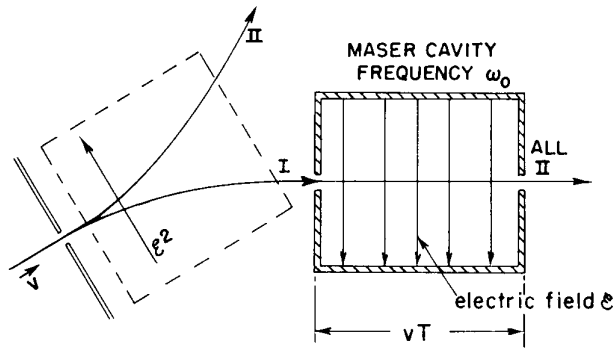


Fig. 9-4. Schematic diagram of the ammonia maser.

9-3 Transitions in a time-dependent field

In the ammonia maser, the beam with molecules in the state $|I\rangle$ and with the energy E_I is sent through a resonant cavity, as shown in Fig. 9-4. The other beam is discarded. Inside the cavity, there will be a time-varying electric field, so the next problem we must discuss is the behavior of a molecule in an electric field that varies with time. We have a completely different kind of a problem—one with a time-varying Hamiltonian. Since H_{ij} depends upon \mathcal{E} , the H_{ij} vary with time, and we must determine the behavior of the system in this circumstance.

To begin with, we write down the equations to be solved:

$$\begin{aligned} i\hbar \frac{dC_1}{dt} &= (E_0 + \mu\mathcal{E})C_1 - AC_2, \\ i\hbar \frac{dC_2}{dt} &= -AC_1 + (E_0 - \mu\mathcal{E})C_2. \end{aligned} \quad (9.36)$$

To be definite, let's suppose that the electric field varies sinusoidally; then we can write

$$\mathcal{E} = 2\mathcal{E}_0 \cos \omega t = \mathcal{E}_0(e^{i\omega t} + e^{-i\omega t}). \quad (9.37)$$

In actual operation the frequency ω will be very nearly equal to the resonant frequency of the molecular transition $\omega_0 = 2A/\hbar$, but for the time being we want to keep things general, so we'll let it have any value at all. The best way to solve our equations is to form linear combinations of C_I and C_{II} as we did before. So we add the two equations, divide by the square root of 2, and use the definitions of C_I and C_{II} that we had in Eq. (9.13). We get

$$i\hbar \frac{dC_{II}}{dt} = (E_0 - A)C_{II} + \mu\mathcal{E}C_I. \quad (9.38)$$

You'll note that this is the same as Eq. (9.9) with an extra term due to the electric field. Similarly, if we subtract the two equations (9.36), we get

$$i\hbar \frac{dC_I}{dt} = (E_0 + A)C_I + \mu\mathcal{E}C_{II}. \quad (9.39)$$

Now the question is, how to solve these equations? They are more difficult than our earlier set, because \mathcal{E} depends on t ; and, in fact, for a general $\mathcal{E}(t)$ the solution is not expressible in elementary functions. However, we can get a good approximation so long as the electric field is small. First we will write

$$\begin{aligned} C_I &= \gamma_I e^{-i(E_0+A)t/\hbar} = \gamma_I e^{-i(E_I)t/\hbar}, \\ C_{II} &= \gamma_{II} e^{-i(E_0-A)t/\hbar} = \gamma_{II} e^{-i(E_{II})t/\hbar}. \end{aligned} \quad (9.40)$$

If there were no electric field, these solutions would be correct with γ_I and γ_{II} just chosen as two complex constants. In fact, since the probability of being in state $|I\rangle$ is the absolute square of C_I and the probability of being in state $|II\rangle$ is the absolute square of C_{II} , the probability of being in state $|I\rangle$ or in state $|II\rangle$ is just $|\gamma_I|^2$ or $|\gamma_{II}|^2$. For instance, if the system were to start originally in state $|II\rangle$ so that γ_I was zero and $|\gamma_{II}|^2$ was one, this condition would go on forever. There would be no chance, if the molecule were originally in state $|II\rangle$, ever to get into state $|I\rangle$.

Now the idea of writing our equations in the form of Eq. (9.40) is that if $\mu\mathcal{E}$ is small in comparison with A , the solutions can still be written in this way, but then γ_I and γ_{II} become slowly varying functions of time—where by “slowly varying” we mean slowly *in comparison* with the exponential functions. That is the trick. We use the fact that γ_I and γ_{II} vary slowly to get an approximate solution.

We want now to substitute C_I from (9.40) in the differential equation (9.39), but we must remember that γ_I is also a function of t . We have

$$i\hbar \frac{dC_I}{dt} = E_I \gamma_I e^{-iE_I t/\hbar} + i\hbar \frac{d\gamma_I}{dt} e^{-iE_I t/\hbar}.$$

The differential equation becomes

$$\left(E_I \gamma_I + i\hbar \frac{d\gamma_I}{dt}\right) e^{-(i/\hbar)E_I t} = E_I \gamma_I e^{-(i/\hbar)E_I t} + \mu\mathcal{E} \gamma_{II} e^{-(i/\hbar)E_{II} t}. \quad (9.41)$$

Similarly, the equation in dC_{II}/dt becomes

$$\left(E_{II} \gamma_{II} + i\hbar \frac{d\gamma_{II}}{dt}\right) e^{-(i/\hbar)E_{II} t} = E_{II} \gamma_{II} e^{-(i/\hbar)E_{II} t} + \mu\mathcal{E} \gamma_I e^{-(i/\hbar)E_I t}. \quad (9.42)$$

Now you will notice that we have equal terms on both sides of each equation. We cancel these terms, and we also multiply the first equation by $e^{+iE_I t/\hbar}$ and the

second by $e^{+iE_{II}t/\hbar}$. Remembering that $(E_I - E_{II}) = 2A = \hbar\omega_0$, we have finally,

$$\begin{aligned} i\hbar \frac{d\gamma_I}{dt} &= \mu\mathcal{E}(t)e^{i\omega_0 t}\gamma_{II}, \\ i\hbar \frac{d\gamma_{II}}{dt} &= \mu\mathcal{E}(t)e^{-i\omega_0 t}\gamma_I. \end{aligned} \quad (9.43)$$

Now we have an apparently simple pair of equations—and they are still exact, of course. The derivative of one variable is a function of time $\mu\mathcal{E}(t)e^{i\omega_0 t}$, multiplied by the second variable; the derivative of the second is a similar time function, multiplied by the first. Although these simple equations cannot be solved in general, we will solve them for some special cases.

We are, for the moment at least, interested only in the case of an oscillating electric field. Taking $\mathcal{E}(t)$ as given in Eq. (9.37), we find that the equations for γ_I and γ_{II} become

$$\begin{aligned} i\hbar \frac{d\gamma_I}{dt} &= \mu\mathcal{E}_0[e^{i(\omega+\omega_0)t} + e^{-i(\omega-\omega_0)t}]\gamma_{II}, \\ i\hbar \frac{d\gamma_{II}}{dt} &= \mu\mathcal{E}_0[e^{i(\omega-\omega_0)t} + e^{-i(\omega+\omega_0)t}]\gamma_I. \end{aligned} \quad (9.44)$$

Now if \mathcal{E}_0 is sufficiently small, the rates of change of γ_I and γ_{II} are also small. The two γ 's will not vary much with t , especially in comparison with the rapid variations due to the exponential terms. These exponential terms have real and imaginary parts that oscillate at the frequency $\omega + \omega_0$ or $\omega - \omega_0$. The terms with $\omega + \omega_0$ oscillate very rapidly about an average value of zero and, therefore, do not contribute very much on the average to the rate of change of γ . So we can make a reasonably good approximation by replacing these terms by their average value, namely, zero. We will just leave them out, and take as our approximation:

$$\begin{aligned} i\hbar \frac{d\gamma_I}{dt} &= \mu\mathcal{E}_0 e^{-i(\omega-\omega_0)t}\gamma_{II}, \\ i\hbar \frac{d\gamma_{II}}{dt} &= \mu\mathcal{E}_0 e^{i(\omega-\omega_0)t}\gamma_I. \end{aligned} \quad (9.45)$$

Even the remaining terms, with exponents proportional to $(\omega - \omega_0)$, will also vary rapidly unless ω is near ω_0 . Only then will the right-hand side vary slowly enough that any appreciable amount will accumulate when we integrate the equations with respect to t . In other words, with a *weak* electric field the only significant frequencies are those near ω_0 .

With the approximation made in getting Eq. (9.45), the equations can be solved exactly, but the work is a little elaborate, so we won't do that until later when we take up another problem of the same type. Now we'll just solve them approximately—or rather, we'll find an exact solution for the case of perfect resonance, $\omega = \omega_0$, and an approximate solution for frequencies near resonance.

9-4 Transitions at resonance

Let's take the case of perfect resonance first. If we take $\omega = \omega_0$, the exponentials are equal to one in both equations of (9.45), and we have just

$$\frac{d\gamma_I}{dt} = -\frac{i\mu\mathcal{E}_0}{\hbar}\gamma_{II}, \quad \frac{d\gamma_{II}}{dt} = -\frac{i\mu\mathcal{E}_0}{\hbar}\gamma_I. \quad (9.46)$$

If we eliminate first γ_I and then γ_{II} from these equations, we find that each satisfies the differential equation of simple harmonic motion:

$$\frac{d^2\gamma}{dt^2} = -\left(\frac{\mu\mathcal{E}_0}{\hbar}\right)^2\gamma. \quad (9.47)$$

The general solutions for these equations can be made up of sines and cosines.

As you can easily verify, the following equations are a solution:

$$\begin{aligned}\gamma_I &= a \cos\left(\frac{\mu\mathcal{E}_0}{\hbar}\right)t + b \sin\left(\frac{\mu\mathcal{E}_0}{\hbar}\right)t, \\ \gamma_{II} &= ib \cos\left(\frac{\mu\mathcal{E}_0}{\hbar}\right)t - ia \sin\left(\frac{\mu\mathcal{E}_0}{\hbar}\right)t,\end{aligned}\tag{9.48}$$

where a and b are constants to be determined to fit any particular physical situation.

For instance, suppose that at $t = 0$ our molecular system was in the upper energy state $|I\rangle$, which would require—from Eq. (9.40)—that $\gamma_I = 1$ and $\gamma_{II} = 0$ at $t = 0$. For this situation we would need $a = 1$ and $b = 0$. The probability that the molecule is in the state $|I\rangle$ at some later t is the absolute square of γ_I , or

$$P_I = |\gamma_I|^2 = \cos^2\left(\frac{\mu\mathcal{E}_0}{\hbar}\right)t.\tag{9.49}$$

Similarly, the probability that the molecule will be in the state $|II\rangle$ is given by the absolute square of γ_{II} ,

$$P_{II} = \gamma_{II}^2 = \sin^2\left(\frac{\mu\mathcal{E}_0}{\hbar}\right)t.\tag{9.50}$$

So long as \mathcal{E} is small and we are on resonance, the probabilities are given by simple oscillating functions. The probability to be in state $|I\rangle$ falls from one to zero and back again, while the probability to be in the state $|II\rangle$ rises from zero to one and back. The time variation of the two probabilities is shown in Fig. 9-5. Needless to say, the sum of the two probabilities is always equal to one; the molecule is always in *some* state!

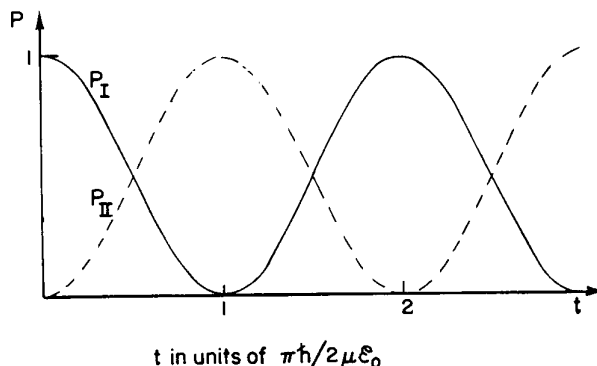


Fig. 9-5. Probabilities for the two states of the ammonia molecule in a sinusoidal electric field.

Let's suppose that it takes the molecule the time T to go through the cavity. If we make the cavity just long enough so that $\mu\mathcal{E}_0T/\hbar = \pi/2$, then a molecule which enters in state $|I\rangle$ will certainly leave it in state $|II\rangle$. If it enters the cavity in the upper state, it will leave the cavity in the lower state. In other words, its energy is decreased, and the loss of energy can't go anywhere else but into the machinery which generates the field. The details by which you can see how the energy of the molecule is fed into the oscillations of the cavity are not simple; however, we don't need to study these details, because we can use the principle of conservation of energy. (We could study them if we had to, but then we would have to deal with the quantum mechanics of the field in the cavity in addition to the quantum mechanics of the atom.)

In summary: the molecule enters the cavity, the cavity field—oscillating at exactly the right frequency—induces transitions from the upper to the lower state, and the energy released is fed into the oscillating field. In an operating maser the molecules deliver enough energy to maintain the cavity oscillations—not only providing enough power to make up for the cavity losses but even providing small amounts of excess power that can be drawn from the cavity. Thus, the molecular energy is converted into the energy of an external electromagnetic field.

Remember that before the beam enters the cavity, we have to use a filter which separates the beam so that only the upper state enters. It is easy to demonstrate that if you were to start with molecules in the lower state, the process will go the other way and take energy out of the cavity. If you put the unfiltered beam in, as many molecules are taking energy out as are putting energy in, so nothing much would happen. In actual operation it isn't necessary, of course, to make $(\mu\epsilon_0 T/\hbar)$ exactly $\pi/2$. For any other value (except an exact integral multiple of π), there is some probability for transitions from state $|I\rangle$ to state $|II\rangle$. For other values, however, the device isn't 100 percent efficient; many of the molecules which leave the cavity could have delivered some energy to the cavity but didn't.

In actual use, the velocity of all the molecules is not the same; they have some kind of Maxwell distribution. This means that the ideal periods of time for different molecules will be different, and it is impossible to get 100 percent efficiency for all the molecules at once. In addition, there is another complication which is easy to take into account, but we don't want to bother with it at this stage. You remember that the electric field in a cavity usually varies from place to place across the cavity. Thus, as the molecules drift across the cavity, the electric field at the molecule varies in a way that is more complicated than the simple sinusoidal oscillation in time that we have assumed. Clearly, one would have to use a more complicated integration to do the problem exactly, but the general idea is still the same.

There are other ways of making masers. Instead of separating the atoms in state $|I\rangle$ from those in state $|II\rangle$ by a Stern-Gerlach apparatus, one can have the atoms already in the cavity (as a gas or a solid) and shift atoms from state $|II\rangle$ to state $|I\rangle$ by some means. One way is one used in the so-called three-state maser. For it, atomic systems are used which have three energy levels, as shown in Fig. 9-6, with the following special properties. The system will absorb radiation (say, light) of frequency $\hbar\omega_1$ and go from the lowest energy level E_{II} to some high-energy level E' , and then will quickly emit photons of frequency $\hbar\omega_2$ and go to the state $|I\rangle$ with energy E_I . The state $|I\rangle$ has a long lifetime so its population can be raised, and the conditions are then appropriate for maser operation between states $|I\rangle$ and $|II\rangle$. Although such a device is called a "three-state" maser, the maser operation really works just as a two-state system such as we are describing.

A laser (Light Amplification by Stimulated Emission of Radiation) is just a maser working at optical frequencies. The "cavity" for a laser usually consists of just two plane mirrors between which standing waves are generated.

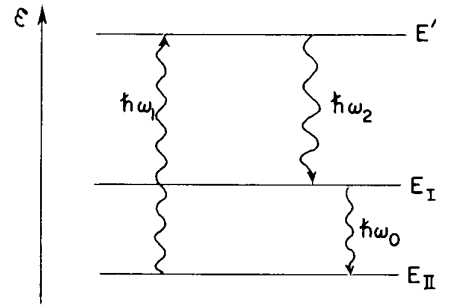


Fig. 9-6. The energy levels of a "three-state" maser.

9-5 Transitions off resonance

Finally, we would like to find out how the states vary in the circumstance that the cavity frequency is nearly, but not exactly, equal to ω_0 . We could solve this problem exactly, but instead of trying to do that, we'll take the important case that the electric field is small and also the period of time T is small, so that $\mu\epsilon_0 T/\hbar$ is much less than one. Then, even in the case of perfect resonance which we have just worked out, the probability of making a transition is small. Suppose that we start again with $\gamma_I = 1$ and $\gamma_{II} = 0$. During the time T we would expect γ_I to remain nearly equal to one, and γ_{II} to remain very small compared with unity. Then the problem is very easy. We can calculate γ_{II} from the second equation in (9.45), taking γ_I equal to one and integrating from $t = 0$ to $t = T$. We get

$$\gamma_{II} = \frac{\mu\epsilon_0}{\hbar} \left[\frac{1 - e^{i(\omega - \omega_0)T}}{\omega - \omega_0} \right]. \quad (9.51)$$

This γ_{II} , used with Eq. (9.40), gives the amplitude to have made a transition from the state $|I\rangle$ to the state $|II\rangle$ during the time interval T . The probability $P(I \rightarrow II)$ to make the transition is $|\gamma_{II}|^2$, or

$$P(I \rightarrow II) = |\gamma_{II}|^2 = \left[\frac{\mu\epsilon_0 T}{\hbar} \right]^2 \frac{\sin^2[(\omega - \omega_0)T/2]}{[(\omega - \omega_0)T/2]^2}. \quad (9.52)$$

It is interesting to plot this probability for a fixed length of time as a function of the frequency of the cavity in order to see how sensitive it is to frequencies near the resonant frequency ω_0 . We show such a plot of $P(I \rightarrow II)$ in Fig. 9-7. (The vertical scale has been adjusted to be 1 at the peak by dividing by the value of the probability when $\omega = \omega_0$.) We have seen a curve like this in the diffraction theory, so you should already be familiar with it. The curve falls rather abruptly to zero for $(\omega - \omega_0) = 2\pi/T$ and never regains significant size for large frequency deviations. In fact, by far the greatest part of the area under the curve lies within the range $\pm\pi/T$. It is possible to show† that the area under the curve is just $2\pi/T$ and is equal to the area of the shaded rectangle drawn in the figure.

Let's examine the implication of our results for a real maser. Suppose that the ammonia molecule is in the cavity for a reasonable length of time, say for one millisecond. Then for $f_0 = 24,000$ megacycles, we can calculate that the probability for a transition falls to zero for a frequency deviation of $(f - f_0)/f_0 = 1/f_0T$, which is five parts in 10^8 . Evidently the frequency must be very close to ω_0 to get a significant transition probability. Such an effect is the basis of the great precision that can be obtained with "atomic" clocks, which work on the maser principle.

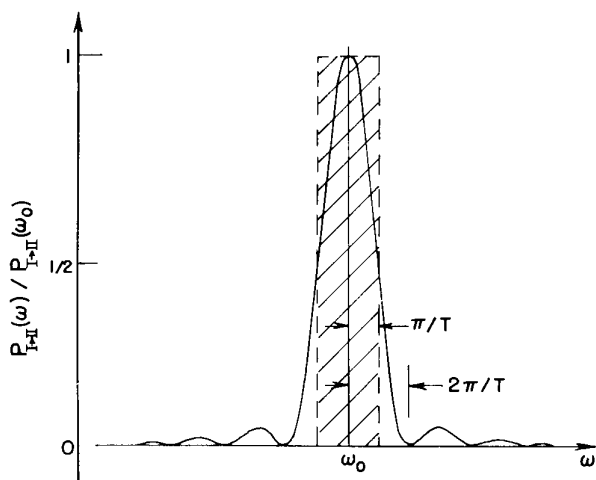


Fig. 9-7. Transition probability for the ammonia molecule as a function of frequency.

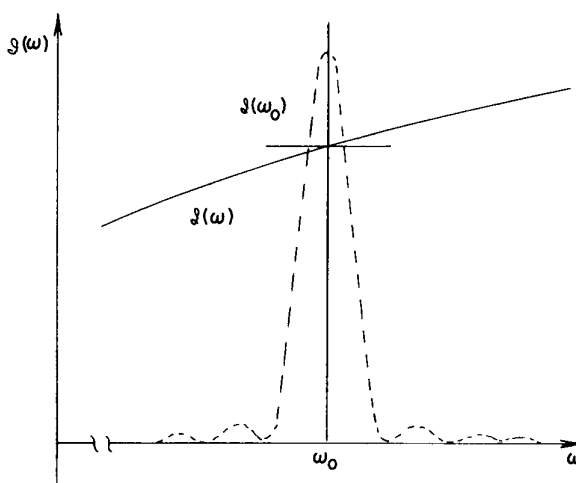


Fig. 9-8. The spectral intensity $g(\omega)$ can be approximated by its value at ω_0 .

9-6 The absorption of light

Our treatment above applies to a more general situation than the ammonia maser. We have treated the behavior of a molecule under the influence of an electric field, whether that field was confined in a cavity or not. So we could be simply shining a beam of "light"—at microwave frequencies—at the molecule and ask for the probability of emission or absorption. Our equations apply equally well to this case, but let's rewrite them in terms of the *intensity* of the radiation rather than the electric field. If we define the intensity g to be the average energy flow per unit area per second, then from Chapter 27 of Volume II, we can write

$$g = \epsilon_0 c^2 |\mathbf{E} \times \mathbf{B}|_{\text{ave}} = \frac{1}{2} \epsilon_0 c^2 (\mathbf{E} \times \mathbf{B})_{\text{max}} = 2 \epsilon_0 c \mathcal{E}_0^2.$$

(The maximum value of \mathcal{E} is $2\mathcal{E}_0$.) The transition probability now becomes:

$$P(I \rightarrow II) = 2\pi \left[\frac{\mu^2}{4\pi\epsilon_0 \hbar^2 c} \right] g T^2 \frac{\sin^2 [(\omega - \omega_0)T/2]}{[(\omega - \omega_0)T/2]^2}. \quad (9.53)$$

† Using the formula $\int_{-\infty}^{\infty} (\sin^2 x/x^2) dx = \pi$.

Ordinarily the light shining on such a system is not exactly monochromatic. It is, therefore, interesting to solve one more problem—that is, to calculate the transition probability when the light has intensity $\mathcal{I}(\omega)$ per unit frequency interval, covering a broad range which includes ω_0 . Then, the probability of going from $|I\rangle$ to $|II\rangle$ will become an integral:

$$P(I \rightarrow II) = 2\pi \left[\frac{\mu^2}{4\pi\epsilon_0\hbar^2c} \right] T^2 \int_0^\infty \mathcal{I}(\omega) \frac{\sin^2[(\omega - \omega_0)T/2]}{[(\omega - \omega_0)T/2]^2} d\omega. \quad (9.54)$$

In general, $\mathcal{I}(\omega)$ will vary much more slowly with ω than the sharp resonance term. The two functions might appear as shown in Fig. 9-8. In such cases, we can replace $\mathcal{I}(\omega)$ by its value $\mathcal{I}(\omega_0)$ at the center of the sharp resonance curve and take it outside of the integral. What remains is just the integral under the curve of Fig. 9-7, which is, as we have seen, just equal to $2\pi/T$. We get the result that

$$P(I \rightarrow II) = 4\pi^2 \left[\frac{\mu^2}{4\pi\epsilon_0\hbar^2c} \right] \mathcal{I}(\omega_0)T. \quad (9.55)$$

This is an important result, because it is the *general theory of the absorption of light by any molecular or atomic system*. Although we began by considering a case in which state $|I\rangle$ had a higher energy than state $|II\rangle$, none of our arguments depended on that fact. Equation (9.55) still holds if the state $|I\rangle$ has a *lower* energy than the state $|II\rangle$; then $P(I \rightarrow II)$ represents the probability for a transition with the *absorption* of energy from the incident electromagnetic wave. The absorption of light by any atomic system always involves the amplitude for a transition in an oscillating electric field between two states separated by an energy $E = \hbar\omega_0$. For any particular case, it is always worked out in just the way we have done here and gives an expression like Eq. (9.55). We, therefore, emphasize the following features of this result. First, the probability is proportional to T . In other words, there is a constant probability per unit time that transitions will occur. Second, this probability is proportional to the *intensity* of the light incident on the system. Finally, the transition probability is proportional to μ^2 , where, you remember, $\mu\mathcal{E}$ defined the shift in energy due to the electric field \mathcal{E} . Because of this, $\mu\mathcal{E}$ also appeared in Eqs. (9.38) and (9.39) as the coupling term that is responsible for the transition between the otherwise stationary states $|I\rangle$ and $|II\rangle$. In other words, for the small \mathcal{E} we have been considering, $\mu\mathcal{E}$ is the so-called “perturbation term” in the Hamiltonian matrix element which connects the states $|I\rangle$ and $|II\rangle$. In the general case, we would have that $\mu\mathcal{E}$ gets replaced by the matrix element $\langle II|H|I\rangle$ (see Section 5-6).

In Volume I (Section 42-5) we talked about the relations among light absorption, induced emission, and spontaneous emission in terms of the Einstein A - and B -coefficients. Here, we have at last the quantum mechanical procedure for computing these coefficients. What we have called $P(I \rightarrow II)$ for our two-state ammonia molecule corresponds precisely to the absorption coefficient B_{nm} of the Einstein radiation theory. For the complicated ammonia molecule—which is too difficult for anyone to calculate—we have taken the matrix element $\langle II|H|I\rangle$ as $\mu\mathcal{E}$, saying that μ is to be gotten from experiment. For simpler atomic systems, the μ_{mn} which belongs to any particular transition can be calculated from the *definition*

$$\mu_{mn}\mathcal{E} = \langle m|H|n\rangle = H_{mn}, \quad (9.56)$$

where H_{mn} is the matrix element of the Hamiltonian which includes the effects of a weak electric field. The μ_{mn} calculated in this way is called the *electric dipole matrix element*. The quantum mechanical theory of the absorption and emission of light is, therefore, reduced to a calculation of these matrix elements for particular atomic systems.

Our study of a simple two-state system has thus led us to an understanding of the general problem of the absorption and emission of light.

Other Two-State Systems

10-1 The hydrogen molecular ion

In the last chapter we discussed some aspects of the ammonia molecule under the approximation that it can be considered as a two-state system. It is, of course, not really a two-state system—there are many states of rotation, vibration, translation, and so on—but each of these states of motion must be analyzed in terms of two internal states because of the flip-flop of the nitrogen atom. Here we are going to consider other examples of systems which, to some approximation or other, can be considered as two-state systems. Lots of things will be approximate because there are always many other states, and in a more accurate analysis they would have to be taken into account. But in each of our examples we will be able to understand a great deal by just thinking about two states.

Since we will only be dealing with two-state systems, the Hamiltonian we need will look just like the one we used in the last chapter. When the Hamiltonian is independent of time, we know that there are two stationary states with definite—and usually different—energies. Generally, however, we start our analysis with a set of base states which are *not* these stationary states, but states which may, perhaps, have some other simple physical meaning. Then, the stationary states of the system will be represented by a linear combination of these base states.

For convenience, we will summarize the important equations from Chapter 9. Let the original choice of base states be $|I\rangle$ and $|2\rangle$. Then any state $|\psi\rangle$ is represented by the linear combination

$$|\psi\rangle = |I\rangle\langle I|\psi\rangle + |2\rangle\langle 2|\psi\rangle = |I\rangle C_1 + |2\rangle C_2. \quad (10.1)$$

The amplitudes C_i (by which we mean either C_1 or C_2) satisfy the two linear differential equations

$$i\hbar \frac{dC_i}{dt} = \sum_j H_{ij} C_j, \quad (10.2)$$

where both i and j take on the values 1 and 2.

When the terms of the Hamiltonian H_{ij} do not depend on t , the two states of definite energy (the stationary states), which we call

$$|\psi_I\rangle = |I\rangle e^{-(i/\hbar)E_I t} \quad \text{and} \quad |\psi_{II}\rangle = |II\rangle e^{-(i/\hbar)E_{II} t},$$

have the energies

$$\begin{aligned} E_I &= \frac{H_{11} + H_{22}}{2} + \sqrt{\left(\frac{H_{11} - H_{22}}{2}\right)^2 + H_{12}H_{21}}, \\ E_{II} &= \frac{H_{11} + H_{22}}{2} - \sqrt{\left(\frac{H_{11} - H_{22}}{2}\right)^2 + H_{12}H_{21}}. \end{aligned} \quad (10.3)$$

The two C 's for each of these states have the same time dependence. The state vectors $|I\rangle$ and $|II\rangle$ which go with the stationary states are related to our original base states $|I\rangle$ and $|2\rangle$ by

$$\begin{aligned} |I\rangle &= |I\rangle a_1 + |2\rangle a_2, \\ |II\rangle &= |I\rangle a'_1 + |2\rangle a'_2. \end{aligned} \quad (10.4)$$

10-1 The hydrogen molecular ion

10-2 Nuclear forces

10-3 The hydrogen molecule

10-4 The benzene molecule

10-5 Dyes

10-6 The Hamiltonian of a spin one-half particle in a magnetic field

10-7 The spinning electron in a magnetic field

The a 's are complex constants, which satisfy

$$|a_1|^2 + |a_2|^2 = 1,$$

$$\frac{a_1}{a_2} = \frac{H_{12}}{E_I - H_{11}}, \quad (10.5)$$

$$|a'_1|^2 + |a'_2|^2 = 1,$$

$$\frac{a'_1}{a'_2} = \frac{H_{12}}{E_{II} - H_{11}}. \quad (10.6)$$

If H_{11} and H_{22} are equal—say both are equal to E_0 —and $H_{12} = H_{21} = -A$, then $E_I = E_0 + A$, $E_{II} = E_0 - A$, and the states $|I\rangle$ and $|II\rangle$ are particularly simple:

$$|I\rangle = \frac{1}{\sqrt{2}} [|1\rangle - |2\rangle], \quad |II\rangle = \frac{1}{\sqrt{2}} [|1\rangle + |2\rangle]. \quad (10.7)$$

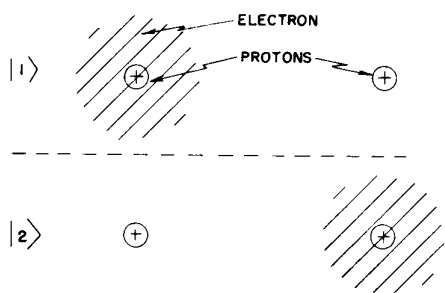


Fig. 10-1. A set of base states for two protons and an electron.

Now we will use these results to discuss a number of interesting examples taken from the fields of chemistry and physics. The first example is the hydrogen molecular ion. A positively ionized hydrogen molecule consists of two protons with one electron worming its way around them. If the two protons are very far apart, what states would we expect for this system? The answer is pretty clear: The electron will stay close to one proton and form a hydrogen atom in its lowest state, and the other proton will remain alone as a positive ion. So, if the two protons are far apart, we can visualize one physical state in which the electron is “attached” to one of the protons. There is, clearly, another state symmetric to that one in which the electron is near the other proton, and the first proton is the one that is an ion. We will take these two as our base states, and we’ll call them $|1\rangle$ and $|2\rangle$. They are sketched in Fig. 10-1. Of course, there are really many states of an electron near a proton, because the combination can exist as any one of the excited states of the hydrogen atom. We are not interested in that variety of states now; we will consider only the situation in which the hydrogen atom is in the lowest state—its ground state—and we will, for the moment, disregard spin of the electron. We can just suppose that for all our states the electron has its spin “up” along the z -axis.[†]

Now to remove an electron from a hydrogen atom requires 13.6 electron volts of energy. So long as the two protons of the hydrogen molecular ion are far apart, it still requires about this much energy—which is for our present considerations a great deal of energy—to get the electron somewhere near the midpoint between the protons. So it is impossible, classically, for the electron to jump from one proton to the other. However, in quantum mechanics it is possible—though not very likely. There is some small amplitude for the electron to move from one proton to the other. As a first approximation, then, each of our base states $|1\rangle$ and $|2\rangle$ will have the energy E_0 , which is just the energy of one hydrogen atom plus one proton. We can take that the Hamiltonian matrix elements H_{11} and H_{22} are both approximately equal to E_0 . The other matrix elements H_{12} and H_{21} , which are the amplitudes for the electron to go back and forth, we will again write as $-A$.

You see that this is the same game we played in the last two chapters. If we disregard the fact that the electron can flip back and forth, we have two states of exactly the same energy. This energy will, however, be split into two energy levels by the possibility of the electron going back and forth—the greater the probability of the transition, the greater the split. So the two energy levels of the system are $E_0 + A$ and $E_0 - A$, and the states which have these definite energies are given by Eqs. (10.7).

[†] This is satisfactory so long as there are no important magnetic fields. We will discuss the effects of magnetic fields on the electron later in this chapter, and the very small effects of spin in the hydrogen atom in Chapter 12.

From our solution we see that if a proton and a hydrogen ion are put anywhere near together, the electron will not stay on one of the protons but will flip back and forth between the two protons. If it starts on one of the protons, it will oscillate back and forth between the states $|I\rangle$ and $|II\rangle$ —giving a time-varying solution. In order to have the lowest energy solution (which does not vary with time), it is necessary to start the system with equal amplitudes for the electron to be around each proton. Remember, there are not two electrons—we are not saying that there is an electron around each proton. There is only *one* electron, and it has the same amplitude— $1/\sqrt{2}$ in magnitude—to be in either position.

Now the amplitude A for an electron which is near one proton to get to the other one depends on the separation between the protons. The closer the protons are together, the larger the amplitude. You remember that we talked in Chapter 7 about the amplitude for an electron to “penetrate a barrier,” which it could not do classically. We have the same situation here. The amplitude for an electron to get across decreases roughly exponentially with the distance—for large distances. Since the transition probability, and therefore A , gets larger when the protons are closer together, the separation of the energy levels will also get larger. If the system is in the state $|I\rangle$, the energy $E_0 + A$ increases with decreasing distance, so these quantum mechanical effects make a *repulsive* force tending to keep the protons apart. On the other hand, if the system is in the state $|II\rangle$, the total energy *decreases* if the protons are brought closer together; there is an *attractive* force pulling the protons together. The variation of the two energies with the distance between the two protons should be roughly as shown in Fig. 10-2. We have, then, a quantum-mechanical explanation of the binding force that holds the H_2^+ ion together.

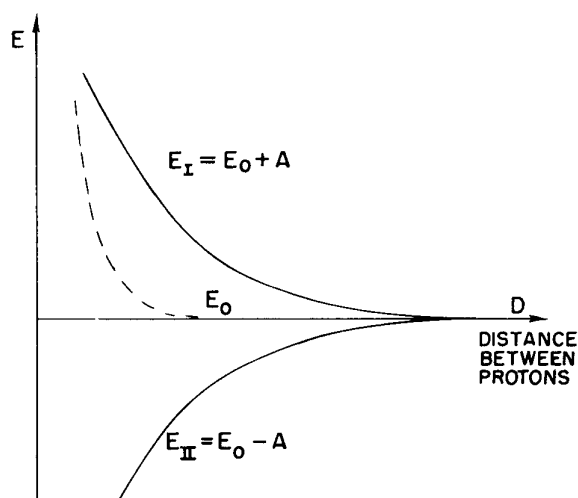


Fig. 10-2. The energies of the two stationary states of the H_2^+ ion as a function of the distance between the two protons.

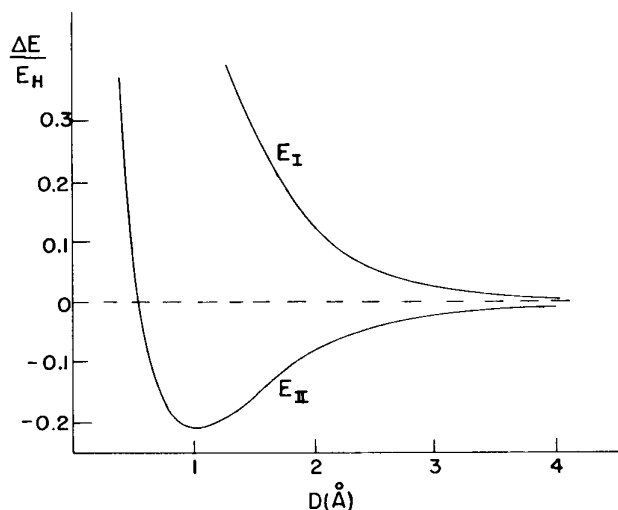


Fig. 10-3. The energy levels of the H_2^+ ion as a function of the interproton distance D . ($E_h = 13.6$ ev.)

We have, however, forgotten one thing. In addition to the force we have just described, there is also an electrostatic repulsive force between the two protons. When the two protons are far apart—as in Fig. 10-1—the “bare” proton sees only a neutral atom, so there is a negligible electrostatic force. At very close distances, however, the “bare” proton begins to get “inside” the electron distribution—that is, it is closer to the proton on the average than to the electron. So there begins to be some extra electrostatic energy which is, of course, positive. This energy—which also varies with the separation—should be included in E_0 . So for E_0 we should take something like the broken-line curve in Fig. 10-2 which rises rapidly for distances less than the radius of a hydrogen atom. We should add and subtract the flip-flop energy A from this E_0 . When we do that, the energies E_I and E_{II} will vary with the interproton distance D as shown in Fig. 10-3. [In this figure, we have plotted the results of a more detailed calculation. The interproton distance

is given in units of $1 \text{ A}(10^{-8} \text{ cm})$, and the excess energy over a proton plus a hydrogen atom is given in units of the binding energy of the hydrogen atom—the so-called “Rydberg” energy, 13.6 ev.] We see that the state $|II\rangle$ has a minimum-energy point. This will be the equilibrium configuration—the lowest energy condition—for the H_2^+ ion. The energy at this point is lower than the energy of a separated proton and hydrogen ion, so the system is bound. A single electron acts to hold the two protons together. A chemist would call it a “one-electron bond.”

This kind of chemical binding is also often called “quantum mechanical resonance” (by analogy with the two coupled pendulums we have described before). But that really sounds more mysterious than it is, it’s only a “resonance” if you start out by making a poor choice for your base states—as we did also! If you picked the state $|II\rangle$, you would have the lowest energy state—that’s all.

We can see in another way why such a state should have a lower energy than a proton and a hydrogen atom. Let’s think about an electron near two protons with some fixed, but not too large, separation. You remember that with a single proton the electron is “spread out” because of the uncertainty principle. It seeks a balance between having a low coulomb *potential* energy and not getting confined into too small a space, which would make a high *kinetic* energy (because of the uncertainty relation $\Delta p \Delta x \approx \hbar$). Now if there are two protons, there is more space where the electron can have a low potential energy. It can spread out—lowering its kinetic energy—without increasing its potential energy. The net result is a lower energy than a hydrogen atom. Then why does the other state $|I\rangle$ have a higher energy? Notice that this state is the *difference* of the states $|1\rangle$ and $|2\rangle$. Because of the symmetry of $|1\rangle$ and $|2\rangle$, the difference must have zero amplitude to find the electron half-way between the two protons. This means that the electron is somewhat more confined, which leads to a larger energy.

We should say that our approximate treatment of the H_2^+ ion as a two-state system breaks down pretty badly once the protons get as close together as they are at the minimum in the curve of Fig. 10-3, and so, will not give a good value for the actual binding energy. For small separations, the energies of the two “states” we imagined in Fig. 6-1 are not really equal to E_0 ; a more refined quantum mechanical treatment is needed.

Suppose we ask now what would happen if instead of two protons, we had two different objects—as, for example, one proton and one lithium positive ion (both particles still with a single positive charge). In such a case, the two terms H_{11} and H_{22} of the Hamiltonian would no longer be equal; they would, in fact, be quite different. If it should happen that the difference ($H_{11} - H_{22}$) is, in absolute value, much greater than $A = -H_{12}$, the attractive force gets very weak, as we can see in the following way.

If we put $H_{12}H_{21} = A^2$ into Eqs. (10.3) we get

$$E = \frac{H_{11} + H_{22}}{2} \pm \frac{H_{11} - H_{22}}{2} \sqrt{1 + \frac{4A^2}{(H_{11} - H_{22})^2}}.$$

When $H_{11} - H_{22}$ is much greater than A^2 , the square root is very nearly equal to

$$1 + \frac{2A^2}{(H_{11} - H_{22})^2}.$$

The two energies are then

$$\begin{aligned} E_I &= H_{11} + \frac{A^2}{(H_{11} - H_{22})}, \\ E_{II} &= H_{22} - \frac{A^2}{(H_{11} - H_{22})}. \end{aligned} \quad (10.8)$$

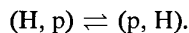
They are now very nearly just the energies H_{11} and H_{22} of the isolated atoms, pushed apart only slightly by the flip-flop amplitude A .

The energy difference $E_I - E_{II}$ is

$$(H_{11} - H_{22}) + \frac{2A^2}{H_{11} - H_{22}}.$$

The *additional* separation from the flip-flop of the electron is no longer equal to $2A$; it is smaller by the factor $A/(H_{11} - H_{22})$, which we are now taking to be much less than one. Also, the dependence of $E_I - E_{II}$ on the separation of the two nuclei is much smaller than for the H_2^+ ion—it is also reduced by the factor $A/(H_{11} - H_{22})$. We can now see why the binding of unsymmetric diatomic molecules is generally very weak.

In our theory of the H_2^+ ion we have discovered an explanation for the mechanism by which an electron shared by two protons provides, in effect, an attractive force between the two protons which can be present even when the protons are at large distances. The attractive force comes from the reduced energy of the system due to the possibility of the electron jumping from one proton to the other. In such a jump the system changes from the configuration (hydrogen atom, proton) to the configuration (proton, hydrogen atom), or switches back. We can write the process symbolically as



The energy shift due to this process is proportional to the amplitude A that an electron whose energy is $-W_H$ (its binding energy in the hydrogen atom) can get from one proton to the other.

For large distances R between the two protons, the electrostatic potential energy of the electron is nearly zero over most of the space it must go when it makes its jump. In this space, then, the electron moves nearly like a free particle in empty space—but with a *negative* energy! We have seen in Chapter 3 [Eq. (3.7)] that the amplitude for a particle of definite energy to get from one place to another a distance r away is proportional to

$$\frac{e^{(i/\hbar)pr}}{r},$$

where p is the momentum corresponding to the definite energy. In the present case (using the nonrelativistic formula), p is given by

$$\frac{p^2}{2m} = -W_H. \quad (10.9)$$

This means that p is an imaginary number,

$$p = i\sqrt{2mW_H}$$

(the other sign for the radical gives nonsense here).

We should expect, then, that the amplitude A for the H_2^+ ion will vary as

$$A \sim \frac{e^{-(\sqrt{2mW_H}/\hbar)R}}{R} \quad (10.10)$$

for large separations R between the two protons. The energy shift due to the electron binding is proportional to A , so there is a force pulling the two protons together which is proportional—for large R —to the derivative of (10.10) with respect to R .

Finally, to be complete, we should remark that in the two-proton, one-electron system there is still one other effect which gives a dependence of the energy on R . We have neglected it until now because it is usually rather unimportant—the exception is just for those very large distances where the energy of the exchange term A has decreased exponentially to very small values. The new effect we are thinking of is the electrostatic attraction of the proton for the hydrogen atom, which comes about in the same way any charged object attracts a neutral object. The bare proton makes an electric field \mathcal{E} (varying as $1/R^2$) at the neutral hydrogen atom. The atom becomes polarized, taking on an induced dipole moment μ proportional to \mathcal{E} . The energy of the dipole is $\mu\mathcal{E}$, which is proportional to \mathcal{E}^2 —or to $1/R^4$. So there is a term in the energy of the system which decreases with the fourth power of the distance. (It is a correction to E_0 .) This energy falls off with

distance more slowly than the shift A given by (10.10); at some large separation R it becomes the only remaining important term giving a variation of energy with R —and, therefore, the only remaining force. Note that the electrostatic term has the same sign for both of the base states (the force is attractive, so the energy is negative) and so also for the two stationary states, whereas the electron exchange term A gives opposite signs for the two stationary states.

10-2 Nuclear forces

We have seen that the system of a hydrogen atom and a proton has an energy of interaction due to the exchange of the single electron which varies at large separations R as

$$\frac{e^{-\alpha R}}{R}, \quad (10.11)$$

with $\alpha = \sqrt{2mW_H}/\hbar$. (One usually says that there is an exchange of a “virtual” electron when—as here—the electron has to jump across a space where it would have a negative energy. More specifically, a “virtual exchange” means that the phenomenon involves a quantum mechanical interference between an exchanged state and a nonexchanged state.)

Now we might ask the following question: Could it be that forces between other kinds of particles have an analogous origin? What about, for example, the nuclear force between a neutron and a proton, or between two protons? In an attempt to explain the nature of nuclear forces, Yukawa proposed that the force between two nucleons is due to a similar exchange effect—only, in this case, due to the virtual exchange, not of an electron, but of a new particle, which he called a “meson.” Today, we would identify Yukawa’s meson with the π -meson (or “pion”) produced in high-energy collisions of protons or other particles.

Let’s see, as an example, what kind of a force we would expect from the exchange of a positive pion (π^+) of mass m_π between a proton and a neutron. Just as a hydrogen atom H^0 can go into a proton p^+ by giving up an electron e^-



a proton p^+ can go into a neutron n^0 by giving up a π^+ meson:



So if we have a proton at a and a neutron at b separated by the distance R , the proton can become a neutron by emitting a π^+ which is then absorbed by the neutron at b , turning it into a proton. There is an energy of interaction of the two-nucleon (plus pion) system which depends on the amplitude A for the pion exchange—just as we found for the electron exchange in the H_2^+ ion.

In the process (10.12), the energy of the H^0 atom is less than that of the proton by W_H (calculating nonrelativistically, and omitting the rest energy mc^2 of the electron), so the electron has a negative *kinetic* energy—or imaginary momentum—as in Eq. (10.9). In the nuclear process (10.13), the proton and neutron have almost equal masses, so the π^+ will have zero *total* energy. The relation between the total energy E and the momentum p for a pion of mass m_π is

$$E^2 = p^2c^2 + m_\pi^2c^4$$

Since E is zero (or at least negligible in comparison with m_π), the momentum is again imaginary:

$$p = im_\pi c.$$

Using the same arguments we gave for the amplitude that a bound electron would penetrate the barrier in the space between two protons, we get for the nuclear case an exchange amplitude A which should—for large R —go as

$$\frac{e^{-(m_\pi c/\hbar)R}}{R}. \quad (10.14)$$

The interaction energy is proportional to A , and so varies in the same way. We get an energy variation in the form of the so-called *Yukawa potential* between two nucleons. Incidentally, we obtained this same formula earlier directly from the differential equation for the motion of a pion in free space [see Chapter 28, Vol. II, Eq. (28.18)].

We can, following the same line of argument, discuss the interaction between two protons (or between two neutrons) which results from the exchange of a *neutral pion* (π^0). The basic process is now

$$p^+ \rightarrow p^+ + \pi^0. \quad (10.15)$$

A proton can emit a virtual π^0 , but then it remains still a proton. If we have two protons, proton No. 1 can emit a virtual π^0 which is absorbed by proton No. 2. At the end, we still have two protons. This is somewhat different from the H_2^+ ion. There the H^0 went into a different condition—the proton—after emitting the electron. Now we are assuming that a proton can emit a π^0 without changing its character. Such processes are, in fact, observed in high-energy collisions. The process is analogous to the way that an electron emits a photon and ends up still an electron:

$$e \rightarrow e + \text{photon}. \quad (10.16)$$

We do not “see” the photons inside the electrons before they are emitted or after they are absorbed, and their emission does not change the “nature” of the electron.

Going back to the two protons, there is an interaction energy which arises from the amplitude A that one proton emits a neutral pion which travels across (with imaginary momentum) to the other proton and is absorbed there. This amplitude is again proportional to (10.14), with m_π the mass of the neutral pion. All the same arguments give an equal interaction energy for two neutrons. Since the nuclear forces (disregarding electrical effects) between neutron and proton, between proton and proton, between neutron and neutron are the same, we conclude that the masses of the charged and neutral pions should be the same. Experimentally, the masses are indeed very nearly equal, and the small difference is about what one would expect from electric self-energy corrections (see Chapter 28, Vol. II).

There are other kinds of particles—like K-mesons—which can be exchanged between two nucleons. It is also possible for two pions to be exchanged at the same time. But all of these other exchanged “objects” have a rest mass m_x higher than the pion mass m_π , and lead to terms in the exchange amplitude which vary as

$$\frac{e^{-(m_x c/\hbar)R}}{R}.$$

These terms die out faster with increasing R than the one-meson term. No one knows, today, how to calculate these higher-mass terms, but for large enough values of R only the one-pion term survives. And, indeed, those experiments which involve nuclear interactions only at large distances do show that the interaction energy is as predicted from the one-pion exchange theory.

In the classical theory of electricity and magnetism, the coulomb electrostatic interaction and the radiation of light by an accelerating charge are closely related—both come out of the Maxwell equations. We have seen in the quantum theory that light can be represented as the quantum excitations of the harmonic oscillations of the classical electromagnetic fields in a box. Alternatively, the quantum theory can be set up by describing light in terms of particles—photons—which obey Bose statistics. We emphasized in Section 4–5 that the two alternative points of view always give identical predictions. Can the second point of view be carried through completely to include *all* electromagnetic effects? In particular, if we want to describe the electromagnetic field purely in terms of Bose particles—that is, in terms of photons—what is the coulomb force due to?

From the “particle” point of view the coulomb interaction between two electrons *comes from the exchange of a virtual photon*. One electron emits a photon—as in reaction (10.16)—which goes over to the second electron, where it is absorbed in the reverse of the same reaction. The interaction energy is again given

by a formula like (10.14), but now with m_π replaced by the rest mass of the photon—which is zero. So the virtual exchange of a photon between two electrons gives an interaction energy that varies simply inversely as R , the distance between the two electrons—just the normal coulomb potential energy! In the “particle” theory of electromagnetism, the process of a virtual photon exchange gives rise to all the phenomena of electrostatics.

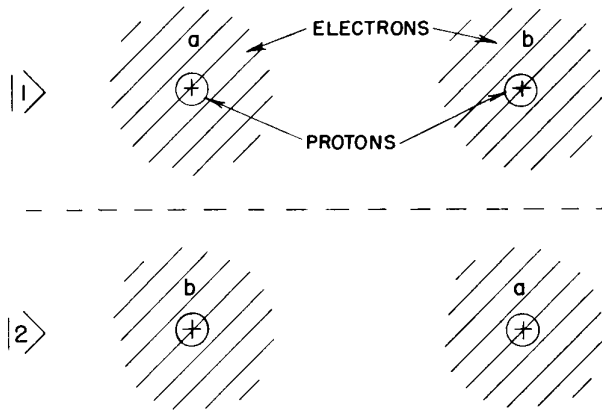


Fig. 10-4. A set of base states for the H_2 molecule.

10-3 The hydrogen molecule

As our next two-state system we will look at the neutral hydrogen molecule H_2 . It is, naturally, more complicated to understand because it has two electrons. Again, we start by thinking of what happens when the two protons are well separated. Only now we have two electrons to add. To keep track of them, we'll call one of them “electron a ” and the other “electron b .” We can again imagine two possible states. One possibility is that “electron a ” is around the first proton and “electron b ” is around the second, as shown in Fig. 10-4(a). We have simply two hydrogen atoms. We will call this state $|1\rangle$. There is also another possibility: that “electron b ” is around the first proton and that “electron a ” is around the second. We call this state $|2\rangle$. From the symmetry of the situation, those two possibilities should be energetically equivalent, but, as we will see, the energy of the system is *not* just the energy of two hydrogen atoms. We should mention that there are many other possibilities. For instance, “electron a ” might be near the first proton and “electron b ” might be in another state around the *same* proton. We'll disregard such a case, since it will certainly have higher energy (because of the large coulomb repulsion between the two electrons). For greater accuracy, we would have to include such states, but we can get the essentials of the molecular binding by considering just the two states of Fig. 10.4. To this approximation we can describe any state by giving the amplitude $\langle 1 | \phi \rangle$ to be in the state $|1\rangle$ and an amplitude $\langle 2 | \phi \rangle$ to be in state $|2\rangle$. In other words, the state vector $|\phi\rangle$ can be written as the linear combination

$$|\phi\rangle = \sum_i |i\rangle \langle i | \phi \rangle.$$

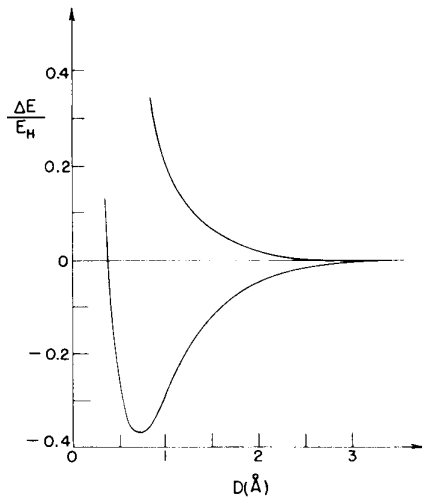


Fig. 10-5. The energy levels of the H_2 molecule for different interproton distances D . ($E_h = 13.6$ ev.)

To proceed, we assume—as usual—that there is some amplitude A that the electrons can move through the intervening space and exchange places. This possibility of exchange means that the energy of the system is split, as we have seen for other two-state systems. As for the hydrogen molecular ion, the splitting is very small when the distance between the protons is large. As the protons approach each other, the amplitude for the electrons to go back and forth increases, so the splitting increases. The decrease of the lower energy state means that there is an attractive force which pulls the atoms together. Again the energy levels rise when the protons get very close together because of the coulomb repulsion. The net final result is that the two stationary states have energies which vary with the separation as shown in Fig. 10-5. At a separation of about 0.74 Å, the lower energy

level reaches a minimum; this is the proton-proton distance of the true hydrogen molecule.

Now you have probably been thinking of an objection. What about the fact that the two electrons are identical particles? We have been calling them “electron *a*” and “electron *b*,” but there really is no way to tell which is which. And we have said in Chapter 4 that for electrons—which are Fermi particles—if there are two ways something can happen by exchanging the electrons, the two amplitudes will interfere with a *negative* sign. This means that if we switch which electron is which, the sign of the amplitude must reverse. We have just concluded, however, that the bound state of the hydrogen molecule would be (at $t = 0$)

$$|II\rangle = \frac{1}{\sqrt{2}}(|I\rangle + |2\rangle).$$

However, according to our rules of Chapter 4, this state is not allowed. If we reverse which electron is which, we get the state

$$\frac{1}{\sqrt{2}}(|2\rangle + |I\rangle),$$

and we get the same sign instead of the opposite one.

These arguments are correct *if both electrons have the same spin*. It is true that if both electrons have spin up (or both have spin down), the only state that is permitted is

$$|I\rangle = \frac{1}{\sqrt{2}}(|I\rangle - |2\rangle).$$

For this state, an interchange of the two electrons gives

$$(|2\rangle - |I\rangle),$$

which is $-|I\rangle$, as required. So if we bring two hydrogen atoms near to each other with their electrons spinning in the same direction, they can go into the state $|I\rangle$ and not state $|II\rangle$. But notice that state $|I\rangle$ is the *upper* energy state. Its curve of energy versus separation has no minimum. The two hydrogens will always repel and will not form a molecule. So we conclude that the hydrogen molecule cannot exist with parallel electron spins. And that is right.

On the other hand, our state $|II\rangle$ is perfectly symmetric for the two electrons. In fact, if we interchange which electron we call *a* and which we call *b* we get back exactly the same state. We saw in Section 4-7 that if two Fermi particles are in the same state, they *must* have opposite spins. So, the bound hydrogen molecule must have one electron with spin up and one with spin down.

The whole story of the hydrogen molecule is really somewhat more complicated if we want to include the proton spins. It is then no longer right to think of the molecule as a *two-state* system. It should really be looked at as an *eight-state* system—there are four possible spin arrangements for each of our states $|I\rangle$ and $|2\rangle$ —so we were cutting things a little short by neglecting the spins. Our final conclusions are, however, correct.

We find that the lowest energy state—the only bound state—of the H_2 molecule has the two electrons with spins opposite. The total spin angular momentum of the electrons is zero. On the other hand, two nearby hydrogen atoms with spins parallel—and so with a total angular momentum \hbar —must be in a higher (unbound) energy state; the atoms repel each other. There is an interesting correlation between the spins and the energies. It gives another illustration of something we mentioned before, which is that there appears to be an “interaction” energy between two spins because the case of parallel spins has a higher energy than the opposite case. In a certain sense you could say that the spins try to reach an antiparallel condition and, in doing so, have the potential to liberate energy—not because there is a large magnetic force, but because of the exclusion principle.

We saw in Section 10-1 that the binding of two *different* ions by a *single* electron is likely to be quite weak. This is *not* true for binding by *two* electrons. Suppose the two protons in Fig. 10-4 were replaced by any two ions (with closed inner electron shells and a single ionic charge), and that the binding energies of an electron at the two ions are different. The energies of states $|1\rangle$ and $|2\rangle$ would still be equal because in each of these states we have one electron bound to each ion. Therefore, we always have the splitting proportional to A . Two-electron binding is ubiquitous—it is the most common valence bond. Chemical binding usually involves this flip-flop game played by two electrons. Although two atoms can be bound together by only one electron, it is relatively rare—because it requires just the right conditions.

Finally, we want to mention that if the energy of attraction for an electron to one nucleus is much greater than to the other, then what we have said earlier about ignoring other possible states is no longer right. Suppose nucleus a (or it may be a positive ion) has a much stronger attraction for an electron than does nucleus b . It may then happen that the total energy is still fairly low even when both electrons are at nucleus a , and no electron is at nucleus b . The strong attraction may more than compensate for the mutual repulsion of the two electrons. If it does, the lowest energy state may have a large amplitude to find both electrons at a (making a negative ion) and a small amplitude to find any electron at b . The state looks like a negative ion with a positive ion. This is, in fact, what happens in an “ionic” molecule like NaCl. You can see that all the gradations between covalent binding and ionic binding are possible.

You can now begin to see how it is that many of the facts of chemistry can be most clearly understood in terms of a quantum mechanical description.

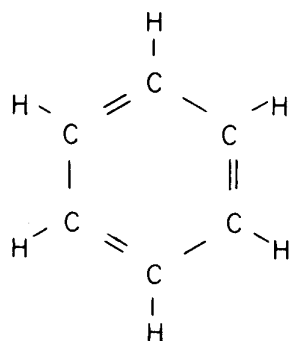


Fig. 10-6. The benzene molecule, C_6H_6 .

10-4 The benzene molecule

Chemists have invented nice diagrams to represent complicated organic molecules. Now we are going to discuss one of the most interesting of them—the benzene molecule shown in Fig. 10-6. It contains six carbon and six hydrogen atoms in a symmetrical array. Each bar of the diagram represents a *pair* of electrons, with spins opposite, doing the covalent bond dance. Each hydrogen atom contributes one electron and each carbon atom contributes four electrons to make up the total of 30 electrons involved. (There are two more electrons close to the nucleus of the carbon which form the first, or K, shell. These are not shown since they are so tightly bound that they are not appreciably involved in the covalent binding.) So each bar in the figure represents a *bond*, or pair of electrons, and the double bonds mean that there are *two pairs* of electrons between alternate pairs of carbon atoms.

There is a mystery about this benzene molecule. We can calculate what energy should be required to form this chemical compound, because the chemists have measured the energies of various compounds which involve pieces of the ring—for instance, they know the energy of a double bond by studying ethylene, and so on. We can, therefore, calculate the total energy we should expect for the benzene

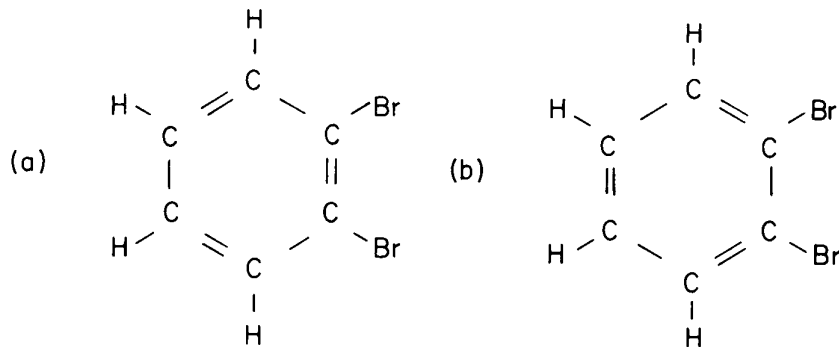


Fig. 10-7. Two possibilities of orthodibromobenzene. The two bromines could be separated by a single bond or by a double bond.

molecule. The actual energy of the benzene ring, however, is much lower than we get by such a calculation; it is more tightly bound than we would expect from what is called an “unsaturated double bond system.” Usually a double bond system which is not in such a ring is easily attacked chemically because it has a relatively high energy—the double bonds can be easily broken by the addition of other hydrogens. But in benzene the ring is quite permanent and hard to break up. In other words, benzene has a much lower energy than you would calculate from the bond picture.

Then there is another mystery. Suppose we replace two adjacent hydrogens by bromine atoms to make ortho-dibromobenzene. There are two ways to do this, as shown in Fig. 10-7. The bromines could be on the opposite ends of a double bond as shown in part (a) of the figure, or could be on the opposite ends of a single bond as in (b). One would think that ortho-dibromobenzene should have two different forms, but it doesn't. There is only one such chemical.†

Now we want to resolve these mysteries—and perhaps you have already guessed how: by noticing, of course, that the “ground state” of the benzene ring is really a two-state system. We could imagine that the bonds in benzene could be in either of the two arrangements shown in Fig. 10-8. You say, “But they are really the same; they should have the same energy.” Indeed, they should. And for that reason they must be analyzed as a two-state system. Each state represents a different configuration of the whole set of electrons, and there is some amplitude A that the whole bunch can switch from one arrangement to the other—there is a chance that the electrons can flip from one dance to the other.

As we have seen, this chance of flipping makes a mixed state whose energy is lower than you would calculate by looking separately at either of the two pictures in Fig. 10-8. Instead, there are two stationary states—one with an energy above and one with an energy below the expected value. So actually, the true normal state (lowest energy) of benzene is neither of the possibilities shown in Fig. 10-8, but it has the amplitude $1/\sqrt{2}$ to be in each of the states shown. It is the only state that is involved in the chemistry of benzene at normal temperatures. Incidentally, the upper state also exists; we can tell it is there because benzene has a strong absorption for ultraviolet light at the frequency $\omega = (E_I - E_{II})/\hbar$. You will remember that in ammonia, where the object flipping back and forth was three protons, the energy separation was in the microwave region. In benzene, the objects are electrons, and because they are much lighter, they find it easier to flip back and forth, which makes the coefficient A very much larger. The result is that the energy difference is much larger—about 1.5 eV, which is the energy of an ultraviolet photon.‡

What happens if we substitute bromine? Again the two “possibilities” (a) and (b) in Fig. 10-7 represent the two different electron configurations. The only difference is that the two base states we start with would have slightly different energies. The lowest energy stationary state will still involve a linear combination of the two states, but with unequal amplitudes. The amplitude for state $|1\rangle$ might have a value something like $\sqrt{2/3}$, say, whereas state $|2\rangle$ would have the magnitude

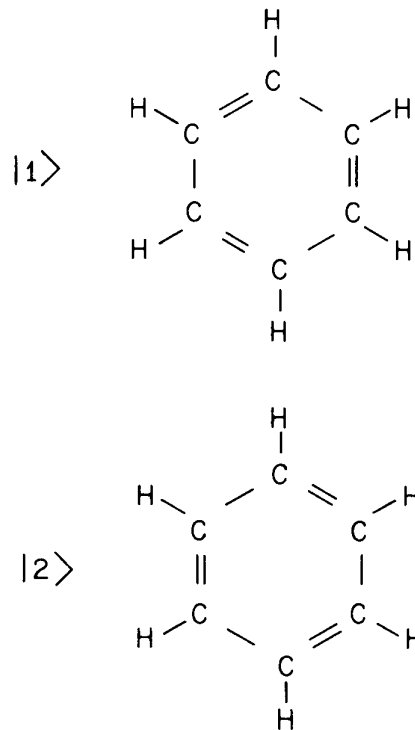


Fig. 10-8. A set of base states for the benzene molecule.

† We are oversimplifying a little. Originally, the chemists thought that there should be *four* forms of dibromobenzene: two forms with the bromines on adjacent carbon atoms (ortho-dibromobenzene), a third form with the bromines on next-nearest carbons (meta-dibromobenzene), and a fourth form with the bromines opposite to each other (para-dibromobenzene). However, they found only three forms—there is only *one* form of the ortho-molecule.

‡ What we have said is a little misleading. Absorption of ultraviolet light would be very weak in the two-state system we have taken for benzene, because the dipole moment matrix element between the two states is zero. [The two states are electrically symmetric, so in our formula Eq. (9.55) for the probability of a transition, the dipole moment μ is zero and no light is absorbed.] If these were the only states, the existence of the upper state would have to be shown in other ways. A more complete theory of benzene, however, which begins with more base states (such as those having adjacent double bonds) shows that the true stationary states of benzene are slightly distorted from the ones we have found. The resulting dipole moments permit the transition we mentioned in the text to occur by the absorption of ultraviolet light.

$\sqrt{1/3}$. We can't say for sure without more information, but once the two energies H_{11} and H_{22} are no longer equal, then the amplitudes C_1 and C_2 no longer have equal magnitudes. This means, of course, that one of the two possibilities in the figure is more likely than the other, but the electrons are mobile enough so that there is some amplitude for both. The other state has different amplitudes (like $\sqrt{1/3}$ and $-\sqrt{2/3}$) but lies at a higher energy. There is only one lowest state, not two as the naive theory of fixed chemical bonds would suggest.

10-5 Dyes

We will give you one more chemical example of the two-state phenomenon—this time on a larger molecular scale. It has to do with the theory of dyes. Many dyes—in fact, most artificial dyes—have an interesting characteristic; they have a kind of symmetry. Figure 10-9 shows an ion of a particular dye called magenta, which has a purplish red color. The molecule has three ring structures—two of which are benzene rings. The third is not exactly the same as a benzene ring because it has only two double bonds inside the ring. The figure shows two equally satisfactory pictures, and we would guess that they should have equal energies. But there is a certain amplitude that all the electrons can flip from one condition to the other, shifting the position of the “unfilled” position to the opposite end. With so many electrons involved, the flipping amplitude is somewhat lower than it is in the case of benzene, and the difference in energy between the two stationary states is smaller. There are, nevertheless, the usual two stationary states $|I\rangle$ and $|II\rangle$ which are the sum and difference combinations of the two base states shown in the figure. The energy separation of $|I\rangle$ and $|II\rangle$ comes out to be equal to the energy of a photon in the optical region. If one shines light on the molecule, there is a very strong absorption at one frequency, and it appears to be brightly colored. That's why it's a dye!

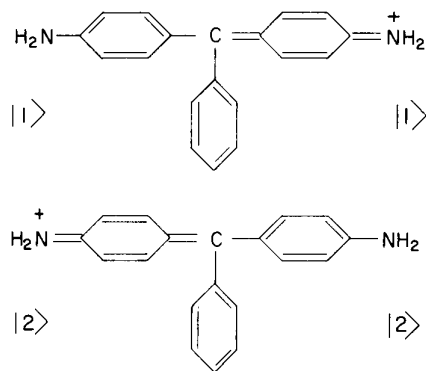


Fig. 10-9. Two base states for the molecule of the dye magenta.

Another interesting feature of such a dye molecule is that in the two base states shown, the center of electric charge is located at different places. As a result, the molecule should be strongly affected by an external electric field. We had a similar effect in the ammonia molecule. Evidently we can analyze it by using exactly the same mathematics, provided we know the numbers E_0 and A . Generally, these are obtained by gathering experimental data. If one makes measurements with many dyes, it is often possible to guess what will happen with some related dye molecule. Because of the large shift in the position of the center of electric charge the value of μ in formula (9.55) is large and the material has a high probability for absorbing light of the characteristic frequency $2A/\hbar$. Therefore, it is not only colored but very strongly so—a small amount of substance absorbs a lot of light.

The rate of flipping—and, therefore, A —is very sensitive to the complete structure of the molecule. By changing A , the energy splitting, and with it the color of the dye, can be changed. Also, the molecules do not have to be perfectly symmetrical. We have seen that the same basic phenomenon exists with slight modifications, even if there is some small asymmetry present. So, one can get some modification of the colors by introducing slight asymmetries in the molecules. For example, another important dye, malachite green, is very similar to magenta, but has two of the hydrogens replaced by CH_3 . It's a different color because the A is shifted and the flip-flop rate is changed.

10-6 The Hamiltonian of a spin one-half particle in a magnetic field

Now we would like to discuss a two-state system involving an object of spin one-half. Some of what we will say has been covered in earlier chapters, but doing it again may help to make some of the puzzling points a little clearer. We can think of an electron at rest as a two-state system. Although we will be talking in this section about “an electron,” what we find out will be true for *any* spin one-half particle. Suppose we choose for our base states $|I\rangle$ and $|2\rangle$ the states in which the z -component of the electron spin is $+\hbar/2$ and $-\hbar/2$.

These states are, of course, the same ones we have called (+) and (−) in earlier chapters. To keep the notation of this chapter consistent, though, we call the “plus” spin state $|1\rangle$ and the “minus” spin state $|2\rangle$ —where “plus” and “minus” refer to the angular momentum in the z -direction.

Any possible state ψ for the electron can be described as in Eq. (10.1) by giving the amplitude C_1 that the electron is in state $|1\rangle$, and the amplitude C_2 that it is in state $|2\rangle$. To treat this problem, we will need to know the Hamiltonian for this two-state system—that is, for an electron in a magnetic field. We begin with the special case of a magnetic field in the z -direction.

Suppose that the vector \mathbf{B} has only a z -component B_z . From the definition of the two base states (that is, spins parallel and antiparallel to \mathbf{B}) we know that they are already stationary states with a definite energy in the magnetic field. State $|1\rangle$ corresponds to an energy† equal to $-\mu B_z$ and state $|2\rangle$ to $+\mu B_z$. The Hamiltonian must be very simple in this case since C_1 , the amplitude to be in state $|1\rangle$, is not affected by C_2 , and vice versa:

$$\begin{aligned} i\hbar \frac{dC_1}{dt} &= E_1 C_1 = -\mu B_z C_1, \\ i\hbar \frac{dC_2}{dt} &= E_2 C_2 = +\mu B_z C_2. \end{aligned} \quad (10.17)$$

For this special case, the Hamiltonian is

$$\begin{aligned} H_{11} &= -\mu B_z, & H_{12} &= 0, \\ H_{21} &= 0, & H_{22} &= +\mu B_z. \end{aligned} \quad (10.18)$$

So we know what the Hamiltonian is for the magnetic field in the z -direction, and we know the energies of the stationary states.

Now suppose the field is *not* in the z -direction. What is the Hamiltonian? How are the matrix elements changed if the field is not in the z -direction? We are going to make an assumption that there is a kind of superposition principle for the terms of the Hamiltonian. More specifically, we want to assume that if two magnetic fields are superposed, the terms in the Hamiltonian simply add—if we know the H_{ij} for a pure B_z and we know the H_{ij} for a pure B_x , then the H_{ij} for both B_z and B_x together is simply the sum. This is certainly true if we consider only fields in the z -direction—if we double B_z , then all the H_{ij} are doubled. So let's assume that H is linear in the field \mathbf{B} . That's all we need to be able to find the H_{ij} for any magnetic field.

Suppose we have a constant field \mathbf{B} . We *could* have chosen our z -axis in its direction, and we *would* have found two stationary states with the energies $\mp\mu B$. Just choosing our axes in a different direction won't change the *physics*. Our *description* of the stationary states will be different, but their energies *will still be* $\mp\mu B$ —that is,

$$\begin{aligned} E_I &= -\mu\sqrt{B_x^2 + B_y^2 + B_z^2} \\ \text{and} \\ E_{II} &= +\mu\sqrt{B_x^2 + B_y^2 + B_z^2}. \end{aligned} \quad (10.19)$$

The rest of the game is easy. We have here the formulas for the energies. We want a Hamiltonian which is linear in B_x , B_y , and B_z , and which will give these energies when used in our general formula of Eq. (10.3). The problem: find the Hamiltonian. First, notice that the energy splitting is symmetric, with an average value of zero. Looking at Eq. (10.3), we can see directly that that requires

$$H_{22} = -H_{11}.$$

(Note that this checks with what we already know when B_x and B_y are both zero;

† We are taking the rest energy $m_0 c^2$ as our “zero” of energy and treating the magnetic moment μ of the electron as a *negative* number, since it points opposite to the spin.

in that case $H_{11} = -\mu B_z$ and $H_{22} = \mu B_z$.) Now if we equate the energies of Eq. (10.3) with what we know from Eq. (10.19), we have

$$\left(\frac{H_{11} - H_{22}}{2}\right)^2 + |H_{12}|^2 = \mu^2(B_x^2 + B_y^2 + B_z^2). \quad (10.20)$$

(We have also made use of the fact that $H_{21} = H_{12}^*$, so that $H_{12}H_{21}$ can also be written as $|H_{12}|^2$.) Again for the special case of a field in the z -direction, this gives

$$\mu^2 B_z^2 + |H_{12}|^2 = \mu^2 B_z^2.$$

Clearly, $|H_{12}|$ must be zero in this special case, which means that H_{12} cannot have any terms in B_z . (Remember, we have said that all terms must be linear in B_x , B_y , and B_z .)

So far, then, we have discovered that H_{11} and H_{22} have terms in B_z , while H_{12} and H_{21} do not. We can make a simple guess that will satisfy Eq. (10.20) if we say that

$$\begin{aligned} H_{11} &= -\mu B_z, \\ H_{22} &= \mu B_z, \end{aligned} \quad (10.21)$$

and

$$|H_{12}|^2 = \mu^2(B_x^2 + B_y^2).$$

And it turns out that that's the *only* way it can be done!

“Wait”—you say—“ H_{12} is not linear in B ; Eq. (10.21) gives $H_{12} = \mu\sqrt{B_x^2 + B_y^2}$.” Not necessarily. There is another possibility which *is* linear, namely,

$$H_{12} = \mu(B_x + iB_y).$$

There are, in fact, several such possibilities—most generally, we could write

$$H_{12} = \mu(B_x \pm iB_y)e^{i\delta},$$

where δ is some arbitrary phase. Which sign and phase should we use? It turns out that you can choose either sign, and any phase you want, and the physical results will always be the same. So the choice is a matter of convention. People ahead of us have chosen to use the minus sign and to take $e^{i\delta} = -1$. We might as well follow suit and write

$$H_{12} = -\mu(B_x - iB_y), \quad H_{21} = -\mu(B_x + iB_y).$$

(Incidentally, these conventions are related to, and consistent with, some of the arbitrary choices we made in Chapter 6.)

The complete Hamiltonian for an electron in an arbitrary magnetic field is, then

$$\begin{aligned} H_{11} &= -\mu B_z, & H_{12} &= -\mu(B_x - iB_y), \\ H_{21} &= -\mu(B_x + iB_y), & H_{22} &= +\mu B_z. \end{aligned} \quad (10.22)$$

And the equations for the amplitudes C_1 and C_2 are

$$\begin{aligned} i\hbar \frac{dC_1}{dt} &= -\mu[B_z C_1 + (B_x - iB_y)C_2], \\ i\hbar \frac{dC_2}{dt} &= -\mu[(B_x + iB_y)C_1 - B_z C_2]. \end{aligned} \quad (10.23)$$

So we have discovered the “equations of motion for the spin states” of an electron in a magnetic field. We guessed at them by making some physical argument, but the real test of any Hamiltonian is that it should give predictions in agreement with experiment. According to any tests that have been made, these equations are right. In fact, although we made our arguments only for constant fields, the Hamiltonian we have written is also right for magnetic fields which vary with time. So we can now use Eq. (10.23) to look at all kinds of interesting problems.

10-7 The spinning electron in a magnetic field

Example number one: We start with a constant field in the z -direction. There are just the two stationary states with energies $\mp\mu B_z$. Suppose we add a small field in the x -direction. Then the equations look like our old two-state problem. We get the flip-flop business once more, and the energy levels are split a little farther apart. Now let's let the x -component of the field vary with time—say, as $\cos \omega t$. The equations are then the same as we had when we put an oscillating electric field on the ammonia molecule in Chapter 9. You can work out the details in the same way. You will get the result that the oscillating field causes transitions from the $+z$ -state to the $-z$ -state—and vice versa—when the horizontal field oscillates near the resonant frequency $\omega_0 = 2\mu B_z/\hbar$. *This gives the quantum mechanical theory of the magnetic resonance phenomena we described in Chapter 35 of Volume II (see Appendix).*

It is also possible to make a maser which uses a spin one-half system. A Stern-Gerlach apparatus is used to produce a beam of particles polarized in, say, the $+z$ -direction, which are sent into a cavity in a constant magnetic field. The oscillating fields in the cavity can couple with the magnetic moment and induce transitions which give energy to the cavity.

Now let's look at the following question. Suppose we have a magnetic field \mathbf{B} which points in the direction whose polar angle is θ and azimuthal angle is ϕ , as in Fig. 10-10. Suppose, additionally, that there is an electron which has been prepared with its spin pointing along this field. What are the amplitudes C_1 and C_2 for such an electron? In other words, calling the state of the electron $|\psi\rangle$, we want to write

$$|\psi\rangle = |1\rangle C_1 + |2\rangle C_2,$$

where C_1 and C_2 are

$$C_1 = \langle 1 | \psi \rangle, \quad C_2 = \langle 2 | \psi \rangle,$$

where by $|1\rangle$ and $|2\rangle$ we mean the same thing we used to call $|+\rangle$ and $|-\rangle$ (referred to our chosen z -axis).

The answer to this question is also in our general equations for two-state systems. First, we know that since the electron's spin is parallel to \mathbf{B} it is in a stationary state with energy $E_I = -\mu B$. Therefore, both C_1 and C_2 must vary as $e^{-iE_I t/\hbar}$, as in (9.18); and their coefficients a_1 and a_2 are given by (10.5), namely,

$$\frac{a_1}{a_2} = \frac{H_{12}}{E_I - H_{11}}. \quad (10.24)$$

An additional condition is that a_1 and a_2 should be normalized so that $|a_1|^2 + |a_2|^2 = 1$. We can take H_{11} and H_{12} from (10.22) using

$$B_z = B \cos \theta, \quad B_x = B \sin \theta \cos \phi, \quad B_y = B \sin \theta \sin \phi.$$

So we have

$$\begin{aligned} H_{11} &= -\mu B \cos \theta, \\ H_{12} &= -\mu B \sin \theta (\cos \phi - i \sin \phi). \end{aligned} \quad (10.25)$$

The last factor in the second equation is, incidentally, $e^{-i\phi}$, so it is simpler to write

$$H_{12} = -\mu B \sin \theta e^{-i\phi}. \quad (10.26)$$

Using these matrix elements in Eq. (10.16)—and canceling $-\mu B$ from numerator and denominator—we find

$$\frac{a_1}{a_2} = \frac{\sin \theta e^{-i\phi}}{1 - \cos \theta}. \quad (10.27)$$

With this ratio and the normalization condition, we can find both a_1 and a_2 . That's not hard, but we can make a short cut with a little trick. Notice that

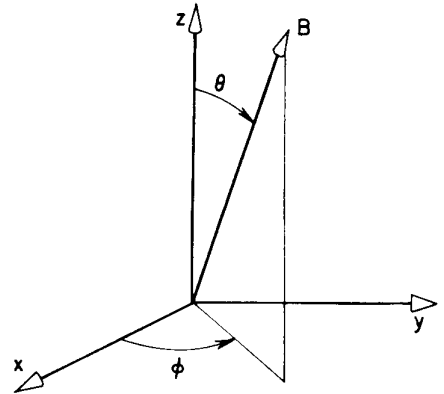


Fig. 10-10. The direction of \mathbf{B} is defined by the polar angle θ and the azimuthal angle ϕ .

$1 - \cos \theta = 2 \sin^2 (\theta/2)$, and that $\sin \theta = 2 \sin (\theta/2) \cos (\theta/2)$. Then Eq. (10.27) is equivalent to

$$\frac{a_1}{a_2} = \frac{\cos \frac{\theta}{2} e^{-i\phi}}{\sin \frac{\theta}{2}}. \quad (10.28)$$

So one possible answer is

$$a_1 = \cos \frac{\theta}{2} e^{-i\phi}, \quad a_2 = \sin \frac{\theta}{2}, \quad (10.29)$$

since it fits with (10.28) and also makes

$$|a_1|^2 + |a_2|^2 = 1.$$

As you know, multiplying both a_1 and a_2 by an arbitrary phase factor doesn't change anything. People generally prefer to make Eqs. (10.29) more symmetric by multiplying both by $e^{i\phi/2}$. So the form usually used is

$$a_1 = \cos \frac{\theta}{2} e^{-i\phi/2}, \quad a_2 = \sin \frac{\theta}{2} e^{+i\phi/2}, \quad (10.30)$$

and this is the answer to our question. The numbers a_1 and a_2 are the amplitudes to find an electron with its spin up or down along the z -axis when we know that its spin is along the axis at θ and ϕ . (The amplitudes C_1 and C_2 are just a_1 and a_2 times $e^{-iE_I t/\hbar}$.)

Now we notice an interesting thing. The strength B of the magnetic field does not appear anywhere in (10.30). The result is clearly the same in the limit that B goes to zero. This means that we have answered *in general* the question of how to represent a particle whose spin is along an arbitrary axis. The amplitudes of (10.30) are the projection amplitudes for spin one-half particles corresponding to the projection amplitudes we gave in Chapter 5 [Eqs. (5.38)] for spin-one particles. We can now find the amplitudes for filtered beams of spin one-half particles to go through any particular Stern-Gerlach filter.

Let $|+z\rangle$ represent a state with spin up along the z -axis, and $|-z\rangle$ represent the spin down state. If $|+z'\rangle$ represents a state with spin up along a z' -axis which makes the polar angles θ and ϕ with the z -axis, then in the notation of Chapter 5, we have

$$\langle +z | +z' \rangle = \cos \frac{\theta}{2} e^{-i\phi/2}, \quad \langle -z | +z' \rangle = \sin \frac{\theta}{2} e^{+i\phi/2}. \quad (10.31)$$

These results are equivalent to what we found in Chapter 6, Eq. (6.36), by purely geometrical arguments. (So if you decided to skip Chapter 6, you now have the essential results anyway.)

As our final example let's look again at one which we've already mentioned a number of times. Suppose that we consider the following problem. We start with an electron whose spin is in some given direction, then turn on a magnetic field in the z -direction for 25 minutes, and then turn it off. What is the final state? Again let's represent the state by the linear combination $|\psi\rangle = |1\rangle C_1 + |2\rangle C_2$. For this problem, however, the states of definite energy are also our base states $|1\rangle$ and $|2\rangle$. So C_1 and C_2 only vary in phase. We know that

$$C_1(t) = C_1(0) e^{-iE_I t/\hbar} = C_1(0) e^{+i\mu B t/\hbar},$$

and

$$C_2(t) = C_2(0) e^{-iE_{II} t/\hbar} = C_2(0) e^{-i\mu B t/\hbar}.$$

Now initially we said the electron spin was set in a given direction. That means that initially C_1 and C_2 are two numbers given by Eqs. (10.30). After we wait for a period of time T , the new C_1 and C_2 are the same two numbers multiplied respectively by $e^{i\mu B_z T/\hbar}$ and $e^{-i\mu B_z T/\hbar}$. What state is that? That's easy. It's exactly the same as if the angle ϕ had been changed by the subtraction of $2\mu B_z T/\hbar$ and the angle θ had been left unchanged. That means that at the end of the time

T , the state $|\psi\rangle$ represents an electron lined up in a direction which differs from the original direction only by a *rotation* about the z -axis through the angle $\Delta\phi = 2\mu B_z T/\hbar$. Since this angle is proportional to T , we can also say the direction of the spin *precesses* at the angular velocity $2\mu B_z/\hbar$ around the z -axis. This result we discussed several times previously in a less complete and rigorous manner. Now we have obtained a complete and accurate quantum mechanical description of the precession of atomic magnets.

It is interesting that the mathematical ideas we have just gone over for the spinning electron in a magnetic field can be applied to *any* two-state system. That means that by making a mathematical *analogy* to the spinning electron, *any problem* about two-state systems can be solved by pure geometry. It works like this. First you shift the zero of energy so that $(H_{11} + H_{22})$ is equal to zero so that $H_{11} = -H_{22}$. Then any two-state problem is *formally* the same as the electron in a magnetic field. All you have to do is *identify* $-\mu B_z$ with H_{11} and $-\mu(B_x - iB_y)$ with H_{12} . No matter what the physics is originally—an ammonia molecule, or whatever—you can translate it into a corresponding electron problem. So if we can solve the electron problem *in general*, we have solved *all* two-state problems.

And we have the general solution for the electron! Suppose you have some state to start with that has spin “up” in some direction, and you have a magnetic field \mathbf{B} that points in some other direction. You just rotate the spin direction around the axis of \mathbf{B} with the *vector* angular velocity $\boldsymbol{\omega}(t)$ equal to a constant times the vector \mathbf{B} (namely $\boldsymbol{\omega} = 2\mu\mathbf{B}/\hbar$). As \mathbf{B} varies with time, you keep moving the axis of the rotation to keep it parallel with \mathbf{B} , and keep changing the speed of rotation so that it is always proportional to the strength of \mathbf{B} . See Fig. 10-11. If you keep doing this, you will end up with a certain final orientation of the spin axis, and the amplitudes C_1 and C_2 are just given by the projections—using (10.30)—into your coordinate frame. You see, it’s just a geometric problem to keep track of where you end up after all the rotating. Although it’s easy to see what’s involved, this geometric problem (of finding the net result of a rotation with a varying angular velocity vector) is not easy to solve explicitly in the general case. Anyway, we see, *in principle*, the general solution to any two-state problem. In the next chapter we will look some more into the mathematical techniques for handling the important case of a spin one-half particle—and, therefore, for handling two-state systems in general.

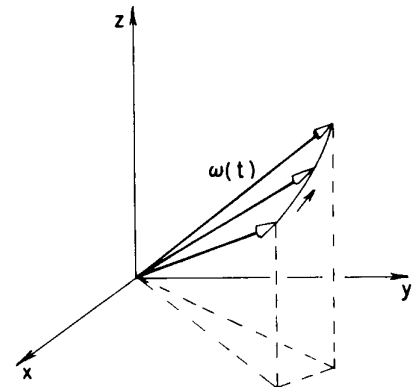


Fig. 10-11. The spin direction of an electron in a varying magnetic field $\mathbf{B}(t)$ precesses at the frequency $\omega(t)$ about an axis parallel to \mathbf{B} .

More Two-State Systems

11-1 The Pauli spin matrices

We continue our discussion of two-state systems. At the end of the last chapter we were talking about a spin one-half particle in a magnetic field. We described the spin state by giving the amplitude C_1 that the z-component of spin angular momentum is $+\hbar/2$ and the amplitude C_2 that it is $-\hbar/2$. In earlier chapters we have called these base states $|+\rangle$ and $|-\rangle$. We will now go back to that notation, although we may occasionally find it convenient to use $|+\rangle$ or $|1\rangle$, and $|-\rangle$ or $|2\rangle$, interchangeably.

We saw in the last chapter that when a spin one-half particle with a magnetic moment μ is in a magnetic field $\mathbf{B} = (B_x, B_y, B_z)$, the amplitudes $C_+(=C_1)$ and $C_-(=C_2)$ are connected by the following differential equations:

$$i\hbar \frac{dC_+}{dt} = -\mu[B_z C_+ + (B_x - iB_y)C_-], \quad (11.1)$$

$$i\hbar \frac{dC_-}{dt} = -\mu[(B_x + iB_y)C_+ - B_z C_-].$$

In other words, the Hamiltonian matrix H_{ij} is

$$\begin{aligned} H_{11} &= -\mu B_z, & H_{12} &= -\mu(B_x - iB_y), \\ H_{21} &= -\mu(B_x + iB_y), & H_{22} &= +\mu B_z. \end{aligned} \quad (11.2)$$

And Eqs. (11.1) are, of course, the same as

$$i\hbar \frac{dC_i}{dt} = \sum_j H_{ij} C_j, \quad (11.3)$$

where i and j take on the values $+$ and $-$ (or 1 and 2).

The two-state system of the electron spin is so important that it is very useful to have a neater way of writing things. We will now make a little mathematical digression to show you how people usually write the equations of a two-state system. It is done this way: First, note that each term in the Hamiltonian is proportional to μ and to some component of \mathbf{B} ; we can then—*purely formally*—write that

$$H_{ij} = -\mu[\sigma_{ij}^x B_x + \sigma_{ij}^y B_y + \sigma_{ij}^z B_z]. \quad (11.4)$$

There is no new physics here; this equation just means that the coefficients σ_{ij}^x , σ_{ij}^y , and σ_{ij}^z —there are $4 \times 3 = 12$ of them—can be figured out so that (11.4) is identical with (11.2).

Let's see what they have to be. We start with B_z . Since B_z appears only in H_{11} and H_{22} , everything will be O.K. if

$$\begin{aligned} \sigma_{11}^z &= 1, & \sigma_{12}^z &= 0, \\ \sigma_{21}^z &= 0, & \sigma_{22}^z &= -1. \end{aligned}$$

We often write the matrix H_{ij} as a little table like this:

$$H_{ij} = \begin{matrix} & \begin{matrix} \xrightarrow{j} \\ \downarrow i \end{matrix} \\ \begin{pmatrix} H_{11} & H_{12} \\ H_{21} & H_{22} \end{pmatrix} \end{matrix}.$$

11-1 The Pauli spin matrices

11-2 The spin matrices as operators

11-3 The solution of the two-state equations

11-4 The polarization states of the photon

11-5 The neutral K-meson†

11-6 Generalization to N -state systems

Review: Chapter 35, Vol. I, *Polarization*

† This section should be omitted on the first reading of this book. It is more advanced than is appropriate in a first course.

For the Hamiltonian of a spin one-half particle in the magnetic field B_z , this is the same as

$$H_{ij} = \overset{i \downarrow}{\underset{j \rightarrow}{\begin{pmatrix} -\mu B_z & -\mu(B_x - iB_y) \\ -\mu(B_x + iB_y) & +\mu B_z \end{pmatrix}}}.$$

In the same way, we can write the coefficients σ_{ij}^z as the matrix

$$\sigma_{ij}^z = \overset{i \downarrow}{\underset{j \rightarrow}{\begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}}}. \quad (11.5)$$

Working with the coefficients of B_x , we get that the terms of σ_x have to be

$$\begin{aligned} \sigma_{11}^x &= 0, & \sigma_{12}^x &= 1, \\ \sigma_{21}^x &= 1, & \sigma_{22}^x &= 0. \end{aligned}$$

Or, in shorthand,

$$\sigma_{ij}^x = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}. \quad (11.6)$$

Finally, looking at B_y , we get

$$\begin{aligned} \sigma_{11}^y &= 0, & \sigma_{12}^y &= -i, \\ \sigma_{21}^y &= i, & \sigma_{22}^y &= 0; \end{aligned}$$

or

$$\sigma_{ij}^y = \begin{pmatrix} 0 & -i \\ i & 0 \end{pmatrix}. \quad (11.7)$$

Table 11-1
The Pauli spin matrices

$\sigma_z = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}$
$\sigma_x = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$
$\sigma_y = \begin{pmatrix} 0 & -i \\ i & 0 \end{pmatrix}$
$1 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$

With these three sigma matrices, Eqs. (11.2) and (11.4) are identical. To leave room for the subscripts i and j , we have shown which σ goes with which component of \mathbf{B} by putting x , y , and z as superscripts. Usually, however, the i and j are omitted—it's easy to imagine they are there—and the x , y , z are written as subscripts. Then Eq. (11.4) is written

$$H = -\mu[\sigma_x B_x + \sigma_y B_y + \sigma_z B_z]. \quad (11.8)$$

Because the sigma matrices are so important—they are used all the time by the professionals—we have gathered them together in Table 11-1. (Anyone who is going to work in quantum physics really has to memorize them.) They are also called the *Pauli spin matrices* after the physicist who invented them.

In the table we have included one more two-by-two matrix which is needed if we want to be able to take care of a system which has two spin states of the same energy, or if we want to choose a different zero energy. For such situations we must add $E_0 C_+$ to the first equation in (11.1) and $E_0 C_-$ to the second equation. We can include this in the new notation if we define the *unit matrix* “1” as δ_{ij} ,

$$1 = \delta_{ij} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, \quad (11.9)$$

and rewrite Eq. (11.8) as

$$H = E_0 \delta_{ij} - \mu(\sigma_x B_x + \sigma_y B_y + \sigma_z B_z). \quad (11.10)$$

Usually, it is *understood* that any constant like E_0 is automatically to be multiplied by the unit matrix; then one writes simply

$$H = E_0 - \mu(\sigma_x B_x + \sigma_y B_y + \sigma_z B_z). \quad (11.11)$$

One reason the spin matrices are useful is that *any* two-by-two matrix at all can be written in terms of them. Any matrix you can write has four numbers in it, say,

$$M = \begin{pmatrix} a & b \\ c & d \end{pmatrix}.$$

It can always be written as a linear combination of four matrices. For example,

$$M = a \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} + b \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix} + c \begin{pmatrix} 0 & 0 \\ 1 & 0 \end{pmatrix} + d \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix}.$$

There are many ways of doing it, but one special way is to say that M is a certain amount of σ_x , plus a certain amount of σ_y , and so on, like this:

$$M = \alpha I + \beta \sigma_x + \gamma \sigma_y + \delta \sigma_z,$$

where the “amounts” α , β , γ , and δ may, in general, be complex numbers.

Since any two-by-two matrix can be represented in terms of the unit matrix and the sigma matrices, we have all that we ever need for *any* two-state system. No matter what the two-state system—the ammonia molecule, the magenta dye, anything—the Hamiltonian equation can be written in terms of the sigmas. Although the sigmas seem to have a geometrical significance in the physical situation of an electron in a magnetic field, they can also be thought of as just useful matrices, which can be used for any two-state problem.

For instance, in one way of looking at things a proton and a neutron can be thought of as the same particle in either of two states. We say the *nucleon* (proton or neutron) is a two-state system—in this case, two states with respect to its charge. When looked at that way, the $|1\rangle$ state can represent the proton and the $|2\rangle$ state can represent the neutron. People say that the nucleon has two “isotopic-spin” states.

Since we will be using the sigma matrices as the “arithmetic” of the quantum mechanics of two-state systems, let’s review quickly the conventions of matrix algebra. By the “sum” of any two or more matrices we mean just what was obvious in Eq. (11.4). In general, if we “add” two matrices A and B , the “sum” C means that each term C_{ij} is given by

$$C_{ij} = A_{ij} + B_{ij}.$$

Each term of C is the sum of the terms in the same slots of A and B .

In Section 5–6 we have already encountered the idea of a matrix “product.” The same idea will be useful in dealing with the sigma matrices. In general, the “product” of two matrices A and B (in that order) is defined to be a matrix C whose elements are

$$C_{ij} = \sum_k A_{ik} B_{kj}. \quad (11.12)$$

It is the sum of products of terms taken in pairs from the i th row of A and the k th column of B . If the matrices are written out in tabular form as in Fig. 11–1, there is a good “system” for getting the terms of the product matrix. Suppose you are calculating C_{23} . You run your left index finger *along* the *second row* of A and your right index finger *down* the *third column* of B , multiplying each pair and adding as you go. We have tried to indicate how to do it in the figure.

$$\begin{pmatrix} A_{11} & A_{12} & A_{13} & A_{14} \\ A_{21} & A_{22} & A_{23} & A_{24} \\ A_{31} & A_{32} & A_{33} & A_{34} \\ A_{41} & A_{42} & A_{43} & A_{44} \end{pmatrix} \cdot \begin{pmatrix} B_{11} & B_{12} & B_{13} & B_{14} \\ B_{21} & B_{22} & B_{23} & B_{24} \\ B_{31} & B_{32} & B_{33} & B_{34} \\ B_{41} & B_{42} & B_{43} & B_{44} \end{pmatrix} = \begin{pmatrix} C_{11} & C_{12} & C_{13} & C_{14} \\ C_{21} & C_{22} & C_{23} & C_{24} \\ C_{31} & C_{32} & C_{33} & C_{34} \\ C_{41} & C_{42} & C_{43} & C_{44} \end{pmatrix}$$

$$C_{ik} = \sum_j A_{ij} B_{jk}$$

$$\text{Example: } C_{23} = A_{21} B_{13} + A_{22} B_{23} + A_{23} B_{33} + A_{24} B_{43}$$

Fig. 11–1. Multiplying two matrices.

It is, of course, particularly simple for two-by-two matrices. For instance, if we multiply σ_x times σ_x , we get

$$\sigma_x^2 = \sigma_x \cdot \sigma_x = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \cdot \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix},$$

which is just the unit matrix 1. Or, for another example, let's work out $\sigma_x \sigma_y$:

$$\sigma_x \sigma_y = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \cdot \begin{pmatrix} 0 & -i \\ i & 0 \end{pmatrix} = \begin{pmatrix} i & 0 \\ 0 & -i \end{pmatrix}.$$

Referring to Table 11-1, you see that the product is just i times the matrix σ_z . (Remember that a number times a matrix just multiplies each term of the matrix.) Since the products of the sigmas taken two at a time are important—as well as rather amusing—we have listed them all in Table 11-2. You can work them out as we have done for σ_x^2 and $\sigma_x \sigma_y$.

There's another very important and interesting point about these σ matrices. We can imagine, if we wish, that the three matrices σ_x , σ_y , and σ_z are analogous to the three components of a vector—it is sometimes called the “sigma vector” and is written σ . It is really a “matrix vector” or a “vector matrix.” It is three different matrices—one matrix associated with each axis, x , y , and z . With it, we can write the Hamiltonian of the system in a nice form which works in any coordinate system:

$$H = -\mu \sigma \cdot B. \quad (11.13)$$

Although we have written our three matrices in the representation in which “up” and “down” are in the z -direction—so that σ_z has a particular simplicity—we could figure out what the matrices would look like in some other representation. Although it takes a lot of algebra, you can show that they change among themselves like the components of a vector. (We won't, however, worry about proving it right now. You can check it if you want.) You can use σ in different coordinate systems as though it is a vector.

You remember that the H is related to energy in quantum mechanics. It is, in fact, just equal to the energy in the simple situation where there is only one state. Even for two-state systems of the electron spin, when we write the Hamiltonian as in Eq. (11.13), it looks very much like the *classical* formula for the energy of a little magnet with magnetic moment μ in a magnetic field. B Classically, we would say

$$U = -\mu \cdot B, \quad (11.14)$$

where μ is the property of the object and B is an external field. We can imagine that Eq. (11.14) can be converted to (11.13) if we replace the classical energy by the Hamiltonian and the classical μ by the matrix $\mu \sigma$. Then, after this purely formal substitution, we interpret the result as a matrix equation. It is sometimes said that to each quantity in classical physics there corresponds a matrix in quantum mechanics. It is really more correct to say that the Hamiltonian matrix corresponds to the energy, and any quantity that can be defined via energy has a corresponding matrix.

For example, the magnetic moment can be defined via energy by saying that the energy in an external field B is $-\mu \cdot B$. This *defines* the magnetic moment vector μ . Then we look at the formula for the Hamiltonian of a real (quantum) object in a magnetic field and try to identify whatever the matrices are that correspond to the various quantities in the classical formula. That's the trick by which *sometimes* classical quantities have their quantum counterparts.

You may try, if you want, to understand how a classical vector is equal to a matrix $\mu \sigma$, and maybe you will discover something—but don't break your head on it. That's not the idea—they are *not equal*. Quantum mechanics is a different kind of a theory to represent the world. It just happens that there are certain correspondences which are hardly more than mnemonic devices—things to remember with. That is, you remember Eq. (11.14) when you learn classical physics;

Table 11-2

Products of the spin matrices

$\sigma_x^2 = 1$
$\sigma_y^2 = 1$
$\sigma_z^2 = 1$
$\sigma_x \sigma_y = -\sigma_y \sigma_x = i \sigma_z$
$\sigma_y \sigma_z = -\sigma_z \sigma_y = i \sigma_x$
$\sigma_z \sigma_x = -\sigma_x \sigma_z = i \sigma_y$

then if you remember the correspondence $\mu \rightarrow \mu\sigma$, you have a handle for remembering Eq. (11.13). Of course, nature knows the quantum mechanics, and the classical mechanics is only an approximation; so there is no mystery in the fact that in classical mechanics there is some shadow of quantum mechanical laws—which are truly the ones underneath. To reconstruct the original object from the shadow is not possible in any direct way, but the shadow does help you to remember what the object looks like. Equation (11.13) is the truth, and Eq. (11.14) is the shadow. Because we learn classical mechanics first, we would like to be able to get the quantum formula from it, but there is no sure-fire scheme for doing that. We must always go back to the real world and discover the correct quantum mechanical equations. When they come out looking like something in classical physics, we are in luck.

If the warnings above seem repetitious and appear to you to be belaboring self-evident truths about the relation of classical physics to quantum physics, please excuse the conditioned reflexes of a professor who has usually taught quantum mechanics to students who hadn't heard about Pauli spin matrices until they were in graduate school. Then they always seemed to be hoping that, somehow, quantum mechanics could be seen to follow as a logical consequence of classical mechanics which they had learned thoroughly years before. (Perhaps they wanted to avoid having to learn something new.) You have learned the classical formula, Eq. (11.14), only a few months ago—and then with warnings that it was inadequate—so maybe you will not be so unwilling to take the quantum formula, Eq. (11.13), as the basic truth.

11-2 The spin matrices as operators

While we are on the subject of mathematical notation, we would like to describe still *another* way of writing things—a way which is used very often because it is so compact. It follows directly from the notation introduced in Chapter 8. If we have a system in a state $|\psi(t)\rangle$, which varies with time, we can—as we did in Eq. (8.31)—write the amplitude that the system would be in the state $|i\rangle$ at $t + \Delta t$ as

$$\langle i | \psi(t + \Delta t) \rangle = \sum_j \langle i | U(t, t + \Delta t) | j \rangle \langle j | \psi(t) \rangle.$$

The matrix element $\langle i | U(t, t + \Delta t) | j \rangle$ is the amplitude that the base state $|j\rangle$ will be converted into the base state $|i\rangle$ in the time interval Δt . We then *defined* H_{ij} by writing

$$\langle i | U(t, t + \Delta t) | j \rangle = \delta_{ij} - \frac{i}{\hbar} H_{ij}(t) \Delta t,$$

and we showed that the amplitudes $C_i(t) = \langle i | \psi(t) \rangle$ were related by the differential equations

$$i\hbar \frac{dC_i}{dt} = \sum_j H_{ij} C_j. \quad (11.15)$$

If we write out the amplitudes C_i explicitly, the same equation appears as

$$i\hbar \frac{d}{dt} \langle i | \psi \rangle = \sum_j H_{ij} \langle j | \psi \rangle. \quad (11.16)$$

Now the matrix elements H_{ij} are also amplitudes which we can write as $\langle i | H | j \rangle$; our differential equation looks like this:

$$i\hbar \frac{d}{dt} \langle i | \psi \rangle = \sum_j \langle i | H | j \rangle \langle j | \psi \rangle. \quad (11.17)$$

We see that $-i/\hbar \langle i | H | j \rangle$ is the amplitude that—under the physical conditions described by H —a state $|j\rangle$ will, during the time dt , “generate” the state $|i\rangle$. (All of this is implicit in the discussion of Section 8-4.)

Now following the ideas of Section 8-2, we can drop out the common term $\langle i |$ in Eq. (11.17)—since it is true for any state $|i\rangle$ —and write that equation simply as

$$i\hbar \frac{d}{dt} |\psi\rangle = \sum_j H |j\rangle \langle j | \psi\rangle. \quad (11.18)$$

Or, going one step further, we can also remove the j and write

$$i\hbar \frac{d}{dt} |\psi\rangle = H |\psi\rangle. \quad (11.19)$$

In Chapter 8 we pointed out that when things are written this way, the H in $H |j\rangle$ or $H |\psi\rangle$ is called an *operator*. From now on we will put the little hat (^) over an operator to remind you that it *is* an operator and not just a number. We will write $\hat{H} |\psi\rangle$. Although the two equations (11.18) and (11.19) *mean exactly the same thing* as Eq. (11.17) or Eq. (11.15), we can *think* about them in a different way. For instance, we would describe Eq. (11.18) in this way: “The time derivative of the *state vector* $|\psi\rangle$ is equal to what you get by operating with the Hamiltonian *operator* \hat{H} on each base state, multiplying by the amplitude $\langle j | \psi\rangle$ that ψ is in the state j , and summing over all j .” Or Eq. (11.19) is described this way. “The time derivative (times $i\hbar$) of a state $|\psi\rangle$ is equal to what you get if you operate with the Hamiltonian \hat{H} on the state vector $|\psi\rangle$.” It’s just a shorthand way of saying what is in Eq. (11.17), but, as you will see, it can be a great convenience.

If we wish, we can carry the “abstraction” idea one more step. Equation (11.19) is true for *any state* $|\psi\rangle$. Also the left-hand side, $i\hbar d/dt$, is also an operator—it’s the operation “differentiate by t and multiply by $i\hbar$.” Therefore, Eq. (11.19) can also be thought of as an equation between operators—the operator equation

$$i\hbar \frac{d}{dt} = \hat{H}.$$

The Hamiltonian operator (within a constant) produces the same result as does d/dt when acting on any state. Remember that this equation—as well as Eq. (11.19)—is *not* a statement that the \hat{H} operator is just the identical *operation* as d/dt . The equations are the dynamical law of nature—the law of motion—for a quantum system.

Just to get some practice with these ideas, we will show you another way we could get to Eq. (11.18). You know that we can write any state $|\psi\rangle$ in terms of its projections into some base set [see Eq. (8.8)],

$$|\psi\rangle = \sum_i |i\rangle \langle i | \psi\rangle. \quad (11.20)$$

How does $|\psi\rangle$ change with time? Well, just take its derivative:

$$\frac{d}{dt} |\psi\rangle = \frac{d}{dt} \sum_i |i\rangle \langle i | \psi\rangle. \quad (11.21)$$

Now the base states $|i\rangle$ do not change with time (at least *we* are always taking them as definite fixed states), but the amplitudes $\langle i | \psi\rangle$ are numbers which may vary. So Eq. (11.21) becomes

$$\frac{d}{dt} |\psi\rangle = \sum_i |i\rangle \frac{d}{dt} \langle i | \psi\rangle. \quad (11.22)$$

Since we know $d\langle i | \psi\rangle/dt$ from Eq. (11.16), we get

$$\begin{aligned} \frac{d}{dt} |\psi\rangle &= -\frac{i}{\hbar} \sum_i |i\rangle \sum_j H_{ij} \langle j | \psi\rangle \\ &= -\frac{i}{\hbar} \sum_{ij} |i\rangle \langle i | H | j\rangle \langle j | \psi\rangle = -\frac{i}{\hbar} \sum_j H |j\rangle \langle j | \psi\rangle. \end{aligned}$$

This is Eq. (11.18) all over again.

So we have many ways of looking at the Hamiltonian. We can think of the set of coefficients H_{ij} as just a bunch of numbers, or we can think of the “amplitudes” $\langle i | H | j \rangle$, or we can think of the “matrix” H_{ij} , or we can think of the operator” \hat{H} . It all means the same thing.

Now let’s go back to our two-state systems. If we write the Hamiltonian in terms of the sigma matrices (with suitable numerical coefficients like B_x , etc.), we can clearly also think of σ_{ij}^x as an amplitude $\langle i | \sigma_x | j \rangle$ or, for short, as the operator $\hat{\sigma}_x$. If we use the operator idea, we can write the equation of motion of a state $|\psi\rangle$ in a magnetic field as

$$i\hbar \frac{d}{dt} |\psi\rangle = -\mu(B_x \hat{\sigma}_x + B_y \hat{\sigma}_y + B_z \hat{\sigma}_z) |\psi\rangle. \quad (11.23)$$

When we want to “use” such an equation we will normally have to express $|\psi\rangle$ in terms of base vectors (just as we have to find the components of space vectors when we want specific numbers). So we will usually want to put Eq. (11.23) in the somewhat expanded form:

$$i\hbar \frac{d}{dt} |\psi\rangle = -\mu \sum_i (B_x \hat{\sigma}_x + B_y \hat{\sigma}_y + B_z \hat{\sigma}_z) |i\rangle \langle i | \psi \rangle. \quad (11.24)$$

Now you will see why the operator idea is so neat. To use Eq. (11.24) we need to know what happens when the $\hat{\sigma}$ operators work on each of the base states. Let’s find out. Suppose we have $\hat{\sigma}_z |+\rangle$; it is some vector $|?\rangle$, but what? Well, let’s multiply it on the left by $\langle + |$; we have

$$\langle + | \hat{\sigma}_z | + \rangle = \sigma_{11} = 1$$

(using Table 11-1). So we know that

$$\langle + | ? \rangle = 1. \quad (11.25)$$

Now let’s multiply $\hat{\sigma}_z |+\rangle$ on the left by $\langle - |$. We get

$$\langle - | \hat{\sigma}_z | + \rangle = \sigma_{21} = 0;$$

so

$$\langle - | ? \rangle = 0. \quad (11.26)$$

There is only one state vector that satisfies both (11.25) and (11.26); it is $|+\rangle$. We discover then that

$$\hat{\sigma}_z |+\rangle = |+\rangle. \quad (11.27)$$

By this kind of argument you can easily show that all of the properties of the sigma matrices can be described in the operator notation by the set of rules given in Table 11-3.

If we have products of sigma matrices, they go over into products of operators. When two operators appear together as a product, you carry out first the operation with the operator which is farthest to the right. For instance, by $\hat{\sigma}_x \hat{\sigma}_y |+\rangle$ we are to understand $\hat{\sigma}_x(\hat{\sigma}_y |+\rangle)$. From Table 11-3, we get $\hat{\sigma}_y |+\rangle = i |-\rangle$, so

$$\hat{\sigma}_x \hat{\sigma}_y |+\rangle = \hat{\sigma}_x(i |-\rangle). \quad (11.28)$$

Now any number—like i —just moves through an operator (operators work only on state vectors); so Eq. (11.28) is the same as

$$\hat{\sigma}_x \hat{\sigma}_y |+\rangle = i \hat{\sigma}_x |-\rangle = i |+\rangle.$$

If you do the same thing for $\hat{\sigma}_x \hat{\sigma}_y |-\rangle$, you will find that

$$\hat{\sigma}_x \hat{\sigma}_y |-\rangle = -i |-\rangle.$$

Looking at Table 11-3, you see that $\hat{\sigma}_x \hat{\sigma}_y$ operating on $|+\rangle$ or $|-\rangle$ gives just what you get if you operate with $\hat{\sigma}_z$ and multiply by $-i$. We can, therefore, say

Table 11-3

Properties of the $\hat{\sigma}$ -operator

$\sigma_z +\rangle = +\rangle$
$\sigma_z -\rangle = - -\rangle$
$\sigma_x +\rangle = -\rangle$
$\sigma_x -\rangle = +\rangle$
$\sigma_y +\rangle = i -\rangle$
$\sigma_y -\rangle = -i +\rangle$

that the operation $\hat{\sigma}_x \hat{\sigma}_y$ is identical with the operation $i\hat{\sigma}_z$, and write this statement as an operator equation:

$$\hat{\sigma}_x \hat{\sigma}_y = i\hat{\sigma}_z. \quad (11.29)$$

Notice that this equation is identical with one of our matrix equations of Table 11-2. So again we see the correspondence between the matrix and operator points of view. Each of the equations in Table 11-2 can, therefore, also be considered as equations about the sigma operators. You can check that they do indeed follow from Table 11-3. It is best, when working with these things, *not* to keep track of whether a quantity like σ or H is an operator or a matrix. All the equations are the same either way, so Table 11-2 is for sigma operators, or for sigma matrices, as you wish.

11-3 The solution of the two-state equations

We can now write our two-state equation in various forms, for example, either as

$$i\hbar \frac{dC_i}{dt} = \sum_j H_{ij} C_j \quad (11.30)$$

or

$$i\hbar \frac{d|\psi\rangle}{dt} = \hat{H}|\psi\rangle.$$

They both mean the same thing. For a spin one-half particle in a magnetic field, the Hamiltonian H is given by Eq. (11.8) or by Eq. (11.13).

If the field is in the z -direction, then—as we have seen several times by now—the solution is that the state $|\psi\rangle$, whatever it is, precesses around the z -axis (just as if you were to take the physical object and rotate it bodily around the z -axis) at an angular velocity equal to twice the magnetic field times μ/\hbar . The same is true, of course, for a magnetic field along any other direction, because the physics is independent of the coordinate system. If we have a situation where the magnetic field varies from time to time in a complicated way, then we can analyze the situation in the following way. Suppose you start with the spin in the $+z$ -direction and you have an x -magnetic field. The spin starts to turn. Then if the x -field is turned off, the spin stops turning. Now if a z -field is turned on, the spin precesses about z , and so on. So depending on how the fields vary in time, you can figure out what the final state is—along what axis it will point. Then you can refer that state back to the original $|+\rangle$ and $|-\rangle$ with respect to z by using the projection formulas we had in Chapter 10 (or Chapter 6). If the state ends up with its spin in the direction (θ, ϕ) , it will have an up-amplitude $\cos(\theta/2)e^{-i\phi/2}$ and a down-amplitude $\sin(\theta/2)e^{+i\phi/2}$. That solves any problem. It is a word description of the solution of the differential equations.

The solution just described is sufficiently general to take care of *any two-state system*. Let's take our example of the ammonia molecule—including the effects of an electric field. If we describe the system in terms of the states $|I\rangle$ and $|II\rangle$, the equations look like this:

$$\begin{aligned} i\hbar \frac{dC_I}{dt} &= +AC_I + \mu\mathcal{E}C_{II}, \\ i\hbar \frac{dC_{II}}{dt} &= -AC_{II} + \mu\mathcal{E}C_I. \end{aligned} \quad (11.31)$$

You say, "No, I remember there was an E_0 in there." Well, we have shifted the origin of energy to make the E_0 zero. (You can always do that by changing both amplitudes by the same factor— $e^{iE_0 T/\hbar}$ —and get rid of any constant energy.) Now if corresponding equations always have the same solutions, then we really don't have to do it twice. If we look at these equations and look at Eq. (11.1), then we can make the following identification. Let's call $|I\rangle$ the state $|+\rangle$ and $|II\rangle$ the state $|-\rangle$. That *does not mean* that we are lining-up the ammonia in space, or that $|+\rangle$ and $|-\rangle$ has anything to do with the z -axis. It is purely artificial.

We have an artificial space that we might “call the ammonia molecule representative space,” or something—a three-dimensional “diagram” in which being “up” corresponds to having the molecule in the state $|I\rangle$ and being “down” along this false z -axis represents having a molecule in the state $|II\rangle$. Then, the equations will be identified as follows. First of all, you see that the Hamiltonian can be written in terms of the sigma matrices as

$$H = +A\sigma_z + \mu\mathcal{E}\sigma_x. \quad (11.32)$$

Or, putting it another way, μB_z in Eq. (11.1) corresponds to $-A$ in Eq. (11.32), and μB_x corresponds to $-\mu\mathcal{E}$. In our “model” space, then, we have a constant B field along the z -direction. If we have an electric field \mathcal{E} which is changing with time, then we have a B field along the x -direction which varies in proportion. *So the behavior of an electron in a magnetic field with a constant component in the z -direction and an oscillating component in the x -direction is mathematically analogous and corresponds exactly to the behavior of an ammonia molecule in an oscillating electric field.* Unfortunately, we do not have the time to go any further into the details of this correspondence, or to work out any of the technical details. We only wished to make the point that *all* systems of two states can be made analogous to a spin one-half object precessing in a magnetic field.

11-4 The polarization states of the photon

There are a number of other two-state systems which are interesting to study, and the first new one we would like to talk about is the photon. To describe a photon we must first give its vector momentum. For a free photon, the frequency is determined by the momentum, so we don’t have to say also what the frequency is. After that, though, we still have a property called the polarization. Imagine that there is a photon coming at you with a definite monochromatic frequency (which will be kept the same throughout all this discussion so that we don’t have a variety of momentum states). Then there are two directions of polarization. In the classical theory, light can be described as having an electric field which oscillates horizontally or an electric field which oscillates vertically (for instance); these two kinds of light are called x -polarized and y -polarized light. The light can also be polarized in some other direction, which can be made up from the superposition of a field in the x -direction and one in the y -direction. Or if you take the x - and the y -components out of phase by 90° , you get an electric field that rotates—the light is elliptically polarized. (This is just a quick reminder of the classical theory of polarized light that we studied in Chapter 35, Vol. I.)

Now, however, suppose we have a *single* photon—just one. There is no electric field that we can discuss in the same way. All we have is *one photon*. But a photon has to have the analog of the classical phenomena of polarization. There must be at least two different kinds of photons. At first, you might think there should be an infinite variety—after all, the electric vector can point in all sorts of directions. We can, however, describe the polarization of a photon as a two-state system. A photon can be in the state $|x\rangle$ or in the state $|y\rangle$. By $|x\rangle$ we mean the polarization state of each one of the photons in a beam of light which *classically* is x -polarized light. On the other hand, by $|y\rangle$ we mean the polarization state of each of the photons in a y -polarized beam. And we can take $|x\rangle$ and $|y\rangle$ as our base states of a photon of given momentum pointing at you—in what we will call the z -direction. So there are two base states $|x\rangle$ and $|y\rangle$, and they are all that are needed to describe any photon at all.

For example, if we have a piece of polaroid set with its axis to pass light polarized in what we call the x -direction, and we send in a photon which we know is in the state $|y\rangle$, it will be absorbed by the polaroid. If we send in a photon which we know is in the state $|x\rangle$, it will come right through as $|x\rangle$. If we take a piece of calcite which takes a beam of polarized light and splits it into an $|x\rangle$ beam and a $|y\rangle$ beam, that piece of calcite is the complete analog of a Stern-Gerlach apparatus which splits a beam of silver atoms into the two states $|+\rangle$ and $|-\rangle$. So every-

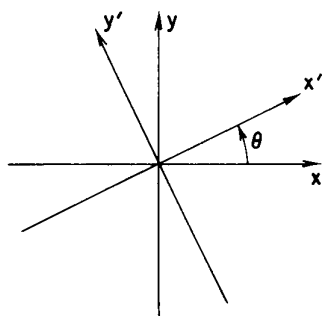


Fig. 11-2. Coordinates at right angles to the momentum vector of the photon.

thing we did before with particles and Stern-Gerlach apparatuses, we can do again with light and pieces of calcite. And what about light filtered through a piece of polaroid set at an angle θ ? Is that another state? Yes, indeed, it *is* another state. Let's call the axis of the polaroid x' to distinguish it from the axes of our base states. See Fig. 11-2. A photon that comes out will be in the state $|x'\rangle$. However, any state can be represented as a linear combination of base states, and the formula for the combination is, here,

$$|x'\rangle = \cos \theta |x\rangle + \sin \theta |y\rangle. \quad (11.33)$$

That is, if a photon comes through a piece of polaroid set at the angle θ (with respect to x), it can still be resolved into $|x\rangle$ and $|y\rangle$ beams—by a piece of calcite, for example. Or you can, if you wish, just analyze it into x - and y -components in your imagination. Either way, you will find the amplitude $\cos \theta$ to be in the $|x\rangle$ state and the amplitude $\sin \theta$ to be in the $|y\rangle$ state.

Now we ask this question: Suppose a photon is polarized in the x' -direction by a piece of polaroid set at the angle θ and arrives at a polaroid at the angle zero—as in Fig. 11-3; what will happen? With what probability will it get through? The answer is the following. After it gets through the first polaroid, it is definitely in the state $|x'\rangle$. The second polaroid will let the photon through if it is in the state $|x\rangle$ (but absorb it if it is the state $|y\rangle$). So we are asking with what probability does the photon appear to be in the state $|x\rangle$? We get that probability from the absolute square of amplitude $\langle x | x' \rangle$ that a photon in the state $|x'\rangle$ is also in the state $|x\rangle$. What is $\langle x | x' \rangle$? Just multiply Eq. (11.33) by $\langle x |$ to get

$$\langle x | x' \rangle = \cos \theta \langle x | x \rangle + \sin \theta \langle x | y \rangle.$$

Now $\langle x | y \rangle = 0$, from the physics—as they *must* be if $|x\rangle$ and $|y\rangle$ are base states—and $\langle x | x \rangle = 1$. So we get

$$\langle x | x' \rangle = \cos \theta,$$

and the probability is $\cos^2 \theta$. For example, if the first polaroid is set at 30° , a photon will get through $3/4$ of the time, and $1/4$ of the time it will heat the polaroid by being absorbed therein.

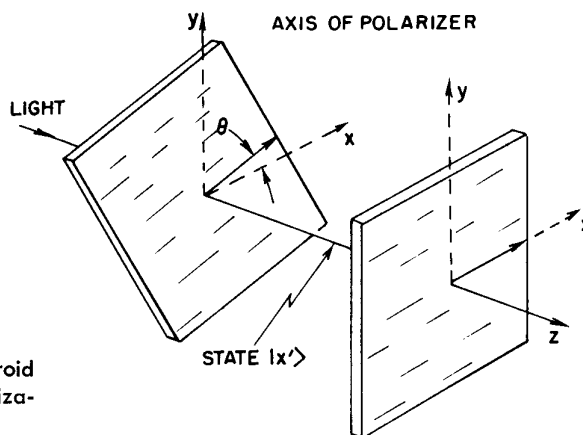


Fig. 11-3. Two sheets of polaroid with angle θ between planes of polarization.

Now let us see what happens classically in the same situation. We would have a beam of light with an electric field which is varying in some way or another—say “unpolarized.” After it gets through the first polaroid, the electric field is oscillating in the x' -direction with a size \mathcal{E} ; we would draw the field as an oscillating vector with a peak value \mathcal{E}_0 in a diagram like Fig. 11-4. Now when the light arrives at the second polaroid, only the x -component, $\mathcal{E}_0 \cos \theta$, of the electric field gets through. The *intensity* is proportional to the square of the field and, therefore, to $\mathcal{E}_0^2 \cos^2 \theta$. So the energy coming through is $\cos^2 \theta$ weaker than the energy which was entering the last polaroid.

The classical picture and the quantum picture give similar results. If you were to throw 10 billion photons at the second polaroid, and the average probability of each one going through is, say, 3/4, you would expect 3/4 of 10 billion would get through. Likewise, the energy that they would carry would be 3/4 of the energy that you attempted to put through. The classical theory says nothing about the statistics of the thing—it simply says that the energy that comes through will be precisely 3/4 of the energy which you were sending in. That is, of course, impossible if there is only one photon. There is no such thing as 3/4 of a photon. It is either *all* there, or it isn't there at all. Quantum mechanics tells us it is *all* there 3/4 of the time. The relation of the two theories is clear.

What about the other kinds of polarization? For example, right-hand circular polarization? In the classical theory, right-hand circular polarization has equal components in x and y which are 90° out of phase. In the quantum theory, a right-hand circularly polarized (RHC) photon has equal amplitudes to be polarized $|x\rangle$ or $|y\rangle$, and the *amplitudes* are 90° out of phase. Calling a RHC photon a state $|R\rangle$ and a LHC photon a state $|L\rangle$, we can write (see Vol. I, Section 33-1)

$$\begin{aligned} |R\rangle &= \frac{1}{\sqrt{2}} (|x\rangle + i|y\rangle), \\ |L\rangle &= -\frac{1}{\sqrt{2}} (|x\rangle - i|y\rangle). \end{aligned} \quad (11.34)$$

—the $1/\sqrt{2}$ is put in to get normalized states. With these states you can calculate any filtering or interference effects you want, using the laws of quantum theory. If you want, you can also choose $|R\rangle$ and $|L\rangle$ as base states and represent everything in terms of them. You only need to show first that $\langle R|L\rangle = 0$ —which you can do by taking the conjugate form of the first equation above [see Eq. (8.13)] and multiplying it by the other. You can resolve light into x - and y -polarizations, or into x' - and y' -polarizations, or into right and left polarizations as a basis.

Just as an example, let's try to turn our formulas around. Can we represent the state $|x\rangle$ as a linear combination of right and left? Yes, here it is:

$$\begin{aligned} |x\rangle &= \frac{1}{\sqrt{2}} (|R\rangle + |L\rangle), \\ |y\rangle &= -\frac{i}{\sqrt{2}} (|R\rangle - |L\rangle). \end{aligned} \quad (11.35)$$

Proof: Add and subtract the two equations in (11.34). It is easy to go from one base to the other.

One curious point has to be made, though. If a photon is right circularly polarized, it shouldn't have anything to do with the x - and y -axes. If we were to look at the same thing from a coordinate system turned at some angle about the direction of flight, the light would still be right circularly polarized—and similarly for left. The right and left circularly polarized light are the same for any such rotation; the definition is independent of any choice of the x -direction (except that the photon direction is given). Isn't that nice—it doesn't take any axes to define it. Much better than x and y . On the other hand, isn't it rather a miracle that when you *add* the right and left together you can find out which direction x was? If "right" and "left" do not depend on x in any way, how is it that we can put them back together again and get x ? We can answer that question in part by writing out the state $|R'\rangle$, which represents a photon RHC polarized in the frame x', y' . In that frame, you would write

$$|R'\rangle = \frac{1}{\sqrt{2}} (|x'\rangle + i|y'\rangle).$$

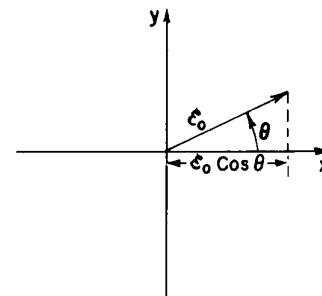


Fig. 11-4. The classical picture of the electric vector \mathcal{E} .

How does such a state look in the frame x, y ? Just substitute x' from Eq.(11. 33) and the corresponding $|y'\rangle$ —we didn't write it down, but it is $(-\sin \theta) |x\rangle + (\cos \theta) |y\rangle$. Then

$$\begin{aligned} |R'\rangle &= \frac{1}{\sqrt{2}} [\cos \theta |x\rangle + \sin \theta |y\rangle - i \sin \theta |x\rangle + i \cos \theta |y\rangle] \\ &= \frac{1}{\sqrt{2}} [(\cos \theta - i \sin \theta) |x\rangle + i(\cos \theta - i \sin \theta) |y\rangle] \\ &= \frac{1}{\sqrt{2}} (|x\rangle + i |y\rangle)(\cos \theta - i \sin \theta). \end{aligned}$$

The first term is just $|R\rangle$, and the second is $e^{+i\theta}$; our result is that

$$|R'\rangle = e^{-i\theta} |R\rangle. \quad (11.36)$$

The states $|R'\rangle$ and $|R\rangle$ are the same except for the phase factor $e^{-i\theta}$. If you work out the same thing for $|L'\rangle$, you get that†

$$|L'\rangle = e^{+i\theta} |L\rangle. \quad (11.37)$$

Now you see what happens. If we add $|R\rangle$ and $|L\rangle$, we get something different from what we get when we add $|R'\rangle$ and $|L'\rangle$. For instance, an x -polarized photon is [Eq. (11.35)] the sum of $|R\rangle$ and $|L\rangle$, but a y -polarized photon is the sum with the phase of one shifted 90° backward and the other 90° forward. That is just what we would get from the sum of $|R'\rangle$ and $|L'\rangle$ for the special angle $\theta = 90^\circ$, and that's right. An x -polarization in the *prime* frame is the same as a y -polarization in the original frame. So it is not exactly true that a circularly polarized photon looks the same for any set of axes. Its *phase* (the phase relation of the right and left circularly polarized states) keeps track of the x -direction.

11-5 The neutral K-meson‡

We will now describe a two-state system in the world of the strange particles—a system for which quantum mechanics gives a most remarkable prediction. To describe it completely would involve us in a lot of stuff about strange particles, so we will, unfortunately, have to cut some corners. We can only give an outline of how a certain discovery was made—to show you the kind of reasoning that was involved. It begins with the discovery by Gell-Mann and Nishijima of the concept of *strangeness* and of a new law of *conservation of strangeness*. It was when Gell-Mann and Pais were analyzing the consequences of these new ideas that they came across the prediction of a most remarkable phenomenon we are going to describe. First, though, we have to tell you a little about “strangeness.”

We must begin with what are called the *strong interactions* of nuclear particles. These are the interactions which are responsible for the strong nuclear forces—as distinct, for instance, from the relatively weaker electromagnetic interactions. The interactions are “strong” in the sense that if two particles get close enough to interact at all, they interact in a big way and produce other particles very easily.

† It's similar to what we found (in Chapter 6) for a spin one-half particle when we rotated the coordinates about the z -axis—then we got the phase factors $e^{\pm i\phi/2}$. It is, in fact, exactly what we wrote down in Section 5-7 for the $|+\rangle$ and $|-\rangle$ states of a spin-one particle—which is no coincidence. The photon is a spin-one particle which has, however, no “zero” state.

‡ We now feel that the material of this section is longer and harder than is appropriate at this point in our development. We suggest that you skip it and continue with Section 11-6. If you are ambitious and have time you may wish to come back to it later. We leave it here, because it is a beautiful example—taken from recent work in high-energy physics—of what can be done with our formulation of the quantum mechanics of two-state systems.

The nuclear particles have also what is called a “weak interaction” by which certain things can happen, such as beta decay, but always very slowly on a nuclear time scale—the weak interactions are many, many orders of magnitude weaker than the strong interactions and even much weaker than electromagnetic interactions.

When the strong interactions were being studied with the big accelerators, people were surprised to find that certain things that “should” happen—that were expected to happen—did not occur. For instance, in some interactions a particle of a certain type did not appear when it was expected. Gell-Mann and Nishijima noticed that many of these peculiar happenings could be explained at once by inventing a new conservation law: the *conservation of strangeness*. They proposed that there was a new kind of attribute associated with each particle—which they called its “strangeness” number—and that in any strong interaction the “quantity of strangeness” is conserved.

Suppose, for instance, that a high-energy negative K-meson—with, say, an energy of many Bev—collides with a proton. Out of the interaction may come many other particles: π -mesons, K-mesons, lambda particles, sigma particles—any of the mesons or baryons listed in Table 2-2 of Vol. I. It is observed, however, that only *certain combinations* appear, and never others. Now certain conservation laws were already known to apply. First, energy and momentum are always conserved. The total energy and momentum after an event must be the same as before the event. Second, there is the conservation of electric charge which says that the total charge of the outgoing particles must be equal to the total charge carried by the original particles. In our example of a K-meson and a proton coming together, the following reactions *do* occur:

$$\begin{aligned} & \text{or} \quad K^- + p \rightarrow p + K^- + \pi^+ + \pi^- + \pi^0 \\ & \quad K^- + p \rightarrow \Sigma^- + \pi^+ \end{aligned} \quad (11.38)$$

We would never get:

$$K^- + p \rightarrow p + K^- + \pi^+ \quad \text{or} \quad K^- + p \rightarrow \Lambda_0 + \pi^+ \quad (11.39)$$

because of the conservation of charge. It was also known that the *number of baryons* is conserved. The number of baryons *out* must be equal to the number of baryons *in*. For this law, an *antiparticle* of a baryon is counted as *minus* one baryon. This means that we can—and *do*—see

$$\begin{aligned} & \text{or} \quad K^- + p \rightarrow \Lambda^0 + \pi^0 \\ & \quad K^- + p \rightarrow p + K^- + p + \bar{p} \end{aligned} \quad (11.40)$$

(where \bar{p} is the antiproton, which carries a negative charge). But we *never* see

$$\begin{aligned} & \text{or} \quad K^- + p \rightarrow K^- + \pi^+ + \pi^0 \\ & \quad K^- + p \rightarrow p + K^- + n \end{aligned} \quad (11.41)$$

(even when there is plenty of energy), because baryons would not be conserved.

These laws, however, do *not* explain the strange fact that the following reactions—which do not immediately appear to be especially different from some of those in (11.38) or (11.40)—are also never observed:

$$\begin{aligned} & \text{or} \quad K^- + p \rightarrow p + K^- + K^0 \\ & \quad K^- + p \rightarrow p + \pi^- \\ & \text{or} \quad K^- + p \rightarrow \Lambda^0 + K^0 \end{aligned} \quad (11.42)$$

The explanation is the conservation of strangeness. With each particle goes a number—its *strangeness* S —and there is a law that in any *strong* interaction, the

Table 11-4

The strangeness numbers of the strongly interacting particles

	S			
	-2	-1	0	+1
Baryons	Ξ^0 Ξ^-	Σ^+ Λ^0, Σ^0 Σ^-	p n	
Mesons		\bar{K}^0 K ⁻	π^+ π^0 π^-	K ⁺ K ⁰

Note: The π^- is the antiparticle of the π^+ (or *vice versa*).

total strangeness *out* must equal the total strangeness that went *in*. The proton and antiproton (p, \bar{p}), the neutron and antineutron (n, \bar{n}), and the π -mesons (π^+ , π^0 , π^-) all have the strangeness number *zero*; the K⁺ and K⁰ mesons have strangeness +1; the K⁻ and \bar{K}^0 (the anti-K⁰),† the Λ^0 and the Σ -particles (+, 0, -) have strangeness -1. There is also a particle with strangeness -2—the Ξ -particle (capital “ksi”)—and perhaps others as yet unknown. We have made a list of these strangenesses in Table 11-4.

Let’s see how the strangeness conservation works in some of the reactions we have written down. If we start with a K⁻ and a proton, we have a total strangeness of $(-1 + 0) = -1$. The conservation of strangeness says that the strangeness of products *after* the reaction must also add up to -1. You see that that is so for the reactions of (11.38) and (11.40). But in the reactions of (11.42) the strangeness of the right-hand side is *zero* in each case. Such reactions do not conserve strangeness, and do not occur. Why? Nobody knows. Nobody knows any more than what we have just told you about this. Nature just works that way.

Now let’s look at the following reaction: a π^- hits a proton. You might, for instance, get a Λ^0 particle plus a neutral K-particle—two neutral particles. Now which neutral K do you get? Since the Λ -particle has a strangeness -1 and the π and p⁺ have a strangeness zero, and since this is a fast production reaction, the strangeness must not change. The K-particle must have strangeness +1—it must therefore be the K⁰. The reaction is

$$\pi^- + p \rightarrow \Lambda^0 + K^0,$$

with

$$S = 0 + 0 = -1 + +1 \quad (\text{conserved}).$$

If the \bar{K}^0 were there instead of the K⁰, the strangeness on the right would be -2—which nature does not permit, since the strangeness on the left side is zero. On the other hand, a \bar{K}^0 can be produced in other reactions, such as

$$n + n \rightarrow n + p + \bar{K}^0 + K^+,$$

$$S = 0 + 0 = 0 + 0 + +1 + -1$$

or

$$K^- + p \rightarrow n + \bar{K}^0,$$

$$S = -1 + 0 = 0 + -1.$$

You may be thinking, “That’s all a lot of stuff, because how do you *know* whether it is a \bar{K}^0 or a K⁰? They look exactly the same. They are antiparticles of each other, so they have exactly the same mass, and both have zero electric charge.

† Read as: “K-naught-bar,” or “K-zero-bar.”

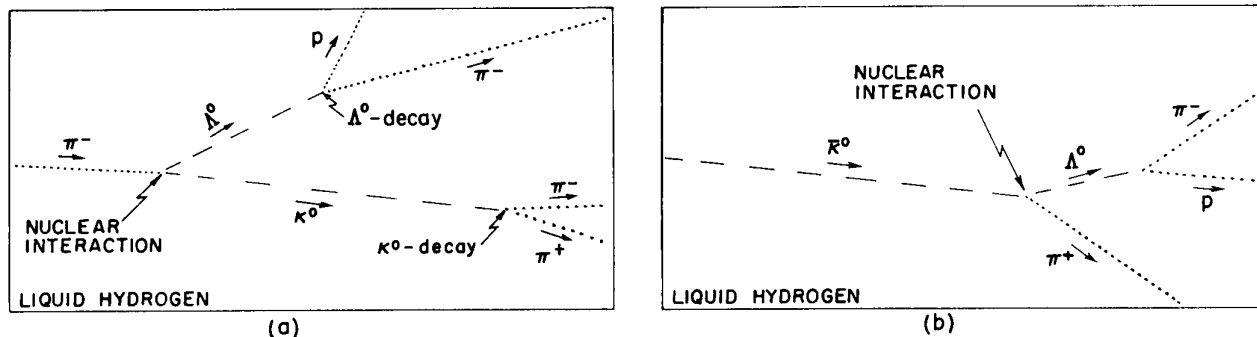
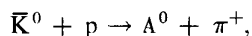


Fig. 11-5. High-energy events as seen in a hydrogen bubble chamber. (a) A π^- meson interacts with a hydrogen nucleus (proton) producing a Λ^0 particle and a K^0 meson. Both particles decay in the chamber. (b) A K^0 meson interacts with a proton producing a π^+ meson and a Λ^0 particle which then decays. (The neutral particles leave no tracks. Their inferred trajectories are indicated here by light dashed lines.)

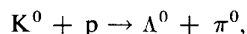
How do you distinguish them?" By the reactions *they* produce. For example, a \bar{K}^0 can interact with matter to produce a Λ -particle, like this:



but a K^0 cannot. There is *no* way a K^0 can produce a Λ -particle when it interacts with ordinary matter (protons and neutrons).† So the experimental distinction between the K^0 and the \bar{K}^0 would be that one of them will and one of them will not produce Λ 's.

One of the predictions of the strangeness theory is then this—if, in an experiment with high-energy pions, a Λ -particle is produced with a neutral K -meson, then *that* neutral K -meson going into other pieces of matter will never produce a Λ . The experiment might run something like this. You send a beam of π^- -mesons into a large hydrogen bubble chamber. A π^- track disappears, but somewhere else a pair of tracks appear (a proton and a π^-) indicating that a Λ -particle has disintegrated‡—see Fig. 11-5. Then you know that there is a K^0 somewhere which you cannot see.

You can, however, figure out where it is going by using the conservation of momentum and energy. [It could reveal itself later by disintegrating into two charged particles, as shown in Fig. 11-5(a).] As the K^0 goes flying along, it may interact with one of the hydrogen nuclei (protons), producing perhaps some other particles. The prediction of the strangeness theory is that it will *never* produce a Λ -particle in a simple reaction like, say,



although a \bar{K}^0 can do just that. That is, in a bubble chamber a \bar{K}^0 might produce the event sketched in Fig. 11-5(b)—in which the Λ^0 is seen because it decays—but a K^0 will not. That's the first part of our story. That's the conservation of strangeness.

The conservation of strangeness is, however, *not perfect*. There are very slow disintegrations of the strange particles—decays taking a long¶ time like 10^{-10} second in which the strangeness is *not* conserved. These are called the "weak" decays. For example, the K^0 disintegrates into a pair of π -mesons (+ and -)

† Except, of course, if it *also* produces *two* K^+ 's or other particles with a total strangeness of +2. We can think here of reactions in which there is insufficient energy to produce these additional strange particles.

‡ The free Λ -particle decays slowly via a *weak* interaction (so strangeness need not be conserved). The decay products are either a p and a π^- , or an n and a π^0 . The lifetime is 2.2×10^{-10} sec.

¶ A typical time for strong interactions is more like 10^{-23} sec.

with a lifetime of 10^{-10} second. That was, in fact, the way K-particles were first seen. Notice that the decay reaction

$$K^0 \rightarrow \pi^+ + \pi^-$$

does not conserve strangeness, so it cannot go “fast” by the strong interaction; it can only go through the weak decay process.

Now the \bar{K}^0 also disintegrates *in the same way*—into a π^+ and a π^- — and also with the same lifetime

$$\bar{K}^0 \rightarrow \pi^- + \pi^+.$$

Again we have a weak decay because it does not conserve strangeness. There is a principle that for any reaction there is the corresponding reaction with “matter” replaced by “antimatter” and *vice versa*. Since the \bar{K}^0 is the antiparticle of the K^0 , it should decay into the antiparticles of the π^+ and π^- , but the antiparticle of a π^+ is the π^- . (Or, if you prefer, *vice versa*. It turns out that for the π -mesons it doesn’t matter which one you call “matter.”) So as a consequence of the weak decays, the K^0 and \bar{K}^0 can go into the same final products. When “seen” through their decays—as in a bubble chamber—they look like the same particle. Only their strong interactions are different.

At last we are ready to describe the work of Gell-Mann and Pais. They first noticed that since the K^0 and the \bar{K}^0 can both turn into states of two π -mesons there must be some amplitude that a K^0 can turn into a \bar{K}^0 , and also that a \bar{K}^0 can turn into a K^0 . Writing the reactions as one does in chemistry, we would have

$$K^0 \rightleftharpoons \pi^- + \pi^+ \rightleftharpoons \bar{K}^0. \quad (11.43)$$

These reactions imply that there is some amplitude per unit time, say $-i/\hbar$ times $\langle \bar{K}^0 | W | K^0 \rangle$, that a K^0 will turn into a \bar{K}^0 through the weak interaction responsible for the decay into two π -mesons. And there is the corresponding amplitude $\langle K^0 | W | \bar{K}^0 \rangle$ for the reverse process. Because matter and antimatter behave in exactly the same way, these two amplitudes are numerically equal; we’ll call them both A :

$$\langle \bar{K}^0 | W | K^0 \rangle = \langle K^0 | W | \bar{K}^0 \rangle = A. \quad (11.44)$$

Now—said Gell-Mann and Pais—here is an interesting situation. What people have been calling two distinct states of the world—the K^0 and the \bar{K}^0 —should really be considered as *one* two-state *system*, because there is an amplitude to go from one state to the other. For a complete treatment, one would, of course, have to deal with more than two states, because there are also the states of 2π ’s, and so on; but since they were mainly interested in the relation of K^0 and \bar{K}^0 , they did not have to complicate things and could make the approximation of a two-state system. The other states *were* taken into account to the extent that their effects appeared implicitly in the amplitudes of Eq. (11.44).

Accordingly, Gell-Mann and Pais analyzed the neutral particle as a two-state system. They began by choosing as their two base states the states $|K^0\rangle$ and $|\bar{K}^0\rangle$. (From here on, the story goes very much as it did for the ammonia molecule.) Any state $|\psi\rangle$ of the neutral K-particle could then be described by giving the amplitudes that it was in either base state. We’ll call these amplitudes

$$C_+ = \langle K^0 | \psi \rangle, \quad C_- = \langle \bar{K}^0 | \psi \rangle. \quad (11.45)$$

The next step was to write the Hamiltonian equations for this two-state system. If there were no coupling between the K^0 and the \bar{K}^0 , the equations would be simply

$$\begin{aligned} i\hbar \frac{dC_+}{dt} &= E_0 C_+, \\ i\hbar \frac{dC_-}{dt} &= E_0 C_-. \end{aligned} \quad (11.46)$$

But since there is the amplitude $\langle \bar{K}^0 | W | K^0 \rangle$ for the K^0 to turn into a \bar{K}^0 there should be the additional term

$$\langle \bar{K}^0 | W | K^0 \rangle C_- = AC_-$$

added to the right-hand side of the first equation. And similarly, the term AC_+ should be inserted in the equation for the rate of change of C_- .

But that's not all. When the two-pion effect is taken into account there is an *additional* amplitude for the K^0 to turn into *itself* through the process

$$K^0 \rightarrow \pi^- + \pi^+ \rightarrow K^0.$$

The additional amplitude, which we would write $\langle K^0 | W | K^0 \rangle$, is just equal to the amplitude $\langle \bar{K}^0 | W | K^0 \rangle$, since the amplitudes to go to and from a pair of π -mesons are identical for the K^0 and the \bar{K}^0 . If you wish, the argument can be written out in detail like this. First write†

$$\langle \bar{K}^0 | W | K^0 \rangle = \langle \bar{K}^0 | W | 2\pi \rangle \langle 2\pi | W | K^0 \rangle$$

and

$$\langle K^0 | W | K^0 \rangle = \langle K^0 | W | 2\pi \rangle \langle 2\pi | W | K^0 \rangle.$$

Because of the symmetry of matter and antimatter

$$\langle 2\pi | W | K^0 \rangle = \langle 2\pi | W | \bar{K}^0 \rangle,$$

and also

$$\langle K^0 | W | 2\pi \rangle = \langle \bar{K}^0 | W | 2\pi \rangle.$$

It then follows that $\langle K^0 | W | K^0 \rangle = \langle \bar{K}^0 | W | K^0 \rangle$, and also that $\langle \bar{K}^0 | W | K^0 \rangle = \langle K^0 | W | \bar{K}^0 \rangle$, as we said earlier. Anyway, there are the two additional amplitudes $\langle K^0 | W | K^0 \rangle$ and $\langle \bar{K}^0 | W | \bar{K}^0 \rangle$, both equal to A , which should be included in the Hamiltonian equations. The first gives a term AC_+ on the right-hand side of the equation for dC_+/dt , and the second gives a new term AC_- in the equation for dC_-/dt . Reasoning this way, Gell-Mann and Pais concluded that the Hamiltonian equations for the $K^0 \bar{K}^0$ system should be

$$\begin{aligned} i\hbar \frac{dC_+}{dt} &= E_0 C_+ + AC_- + AC_+, \\ i\hbar \frac{dC_-}{dt} &= E_0 C_- + AC_+ + AC_-. \end{aligned} \tag{11.47}$$

We must now correct something we have said in earlier chapters: that two amplitudes like $\langle K^0 | W | \bar{K}^0 \rangle$ and $\langle \bar{K}^0 | W | K^0 \rangle$ which are the reverse of each other, are always complex conjugates. That was true when we were talking about particles that did not decay. But if particles can decay—and can, therefore, become “lost”—the two amplitudes are not necessarily complex conjugates. So the equality of (11.44) does not mean that the amplitudes are real numbers; they are in fact complex numbers. The coefficient A is, therefore, complex; and we can't just incorporate it into the energy E_0 .

Having played often with electron spins and such, our heroes knew that the Hamiltonian equations of (11.47) meant that there was *another* pair of base states which could also be used to represent the K -particle system and which would have especially simple behaviors. They said, “Let's take the sum and difference of these two equations. Also, let's measure all our energies from E_0 , and use units for

† We are making a simplification here. The 2π -system can have many states corresponding to various momenta of the π -mesons, and we should make the right-hand side of this equation into a sum over the various base states of the π 's. The complete treatment still leads to the same conclusions.

energy and time that make $\hbar = 1$." (That's what modern theoretical physicists always do. It doesn't change the physics but makes the equations take on a simple form.) Their result:

$$i \frac{d}{dt} (C_+ + C_-) = 2A(C_+ + C_-), \quad i \frac{d}{dt} (C_+ - C_-) = 0. \quad (11.48)$$

It is apparent that the combinations of amplitudes $(C_+ + C_-)$ and $(C_+ - C_-)$ act independently from each other (corresponding, of course, to the stationary states we have been studying earlier). So they concluded that it would be more convenient to use a different representation for the K-particle. They defined the two states

$$|K_1\rangle = \frac{1}{\sqrt{2}} (|K^0\rangle + |\bar{K}^0\rangle), \quad |K_2\rangle = \frac{1}{\sqrt{2}} (|K^0\rangle - |\bar{K}^0\rangle). \quad (11.49)$$

They said that instead of thinking of the K^0 and \bar{K}^0 mesons, we can equally well think in terms of the two "particles" (that is, "states") K_1 and K_2 . (These correspond, of course, to the states we have usually called $|I\rangle$ and $|II\rangle$). We are not using our old notation because we want now to follow the notation of the original authors—and the one you will see in physics seminars.)

Now Gell-Mann and Pais didn't do all this just to get different names for the particles—there is also some strange new physics in it. Suppose that C_1 and C_2 are the amplitudes that some state $|\psi\rangle$ will be either a K_1 or a K_2 meson:

$$C_1 = \langle K_1 | \psi \rangle, \quad C_2 = \langle K_2 | \psi \rangle.$$

From the equations of (11.49),

$$C_1 = \frac{1}{\sqrt{2}} (C_+ + C_-), \quad C_2 = \frac{1}{\sqrt{2}} (C_+ - C_-). \quad (11.50)$$

Then the Eqs. (11.48) become

$$i \frac{dC_1}{dt} = 2AC_1, \quad i \frac{dC_2}{dt} = 0. \quad (11.51)$$

The solutions are

$$C_1(t) = C_1(0)e^{-i2At}, \quad C_2(t) = C_2(0), \quad (11.52)$$

where, of course, $C_1(0)$ and $C_2(0)$ are the amplitudes at $t = 0$.

These equations say that if a neutral K-particle starts out in the state $|K_1\rangle$ at $t = 0$ [then $C_1(0) = 1$ and $C_2(0) = 0$], the amplitudes at the time t are

$$C_1(t) = e^{-i2At}, \quad C_2(t) = 0.$$

Remembering that A is a complex number, it is convenient to take $A = \alpha - i\beta$. (Since the imaginary part of $2A$ turns out to be negative, we write it as *minus* $i\beta$.) With this substitution, $C_1(t)$ reads

$$C_1(t) = C_1(0)e^{-\beta t}e^{-i\alpha t}. \quad (11.53)$$

The probability of finding a K_1 particle at t is the absolute square of this amplitude, which is $e^{-2\beta t}$. And, from Eqs. (11.52), the probability of finding the K_2 state at any time is zero. That means that if you make a K-particle in the state $|K_1\rangle$, the probability of finding it in the same state decreases exponentially with time—but you will never find it in state $|K_2\rangle$. Where does it go? It disintegrates into two π -mesons with the mean life $\tau = 1/2\beta$ which is, experimentally, 10^{-10} sec. We made provisions for that when we said that A was complex.

On the other hand, Eq. (11.52) says that if we make a K-particle completely in the K_2 state, it stays that way forever. Well, that's not really true. It is observed experimentally to disintegrate into *three* π -mesons, but 600 times slower than the

two-pion decay we have described. So there are some other small terms we have left out in our approximation. But so long as we are considering only the two-pion decay, the K_2 lasts “forever.”

Now to finish the story of Gell-Mann and Pais. They went on to consider what happens when a K-particle is produced *with a Λ^0 particle* in a *strong* interaction. Since it must then have a strangeness of +1, it must be produced in the K^0 state. So at $t = 0$ it is neither a K_1 nor a K_2 but a *mixture*. The initial conditions are

$$C_+(0) = 1, \quad C_-(0) = 0.$$

But that means—from Eq. (11.50)—that

$$C_1(0) = \frac{1}{\sqrt{2}}, \quad C_2(0) = \frac{1}{\sqrt{2}},$$

and—from Eq. (11.51)—that

$$C_1(t) = \frac{1}{\sqrt{2}} e^{-\beta t} e^{-i\alpha t}, \quad C_2(t) = \frac{1}{\sqrt{2}}. \quad (11.54)$$

Now remember that K_1 and K_2 are each linear combinations of K^0 and \bar{K}^0 . In Eqs. (11.54) the amplitudes have been chosen so that at $t = 0$ the \bar{K}^0 parts cancel each other out by interference, leaving only a K^0 state. But the $|K_1\rangle$ state *changes with time*, and the $|K_2\rangle$ state *does not*. After $t = 0$ the interference of C_1 and C_2 will give finite amplitudes for both K^0 and \bar{K}^0 .

What does all this mean? Let’s go back and think of the experiment we sketched in Fig. 11-5. A π^- meson has produced a Λ^0 particle and a K^0 meson which is tooting along through the hydrogen in the chamber. As it goes along, there is some small but uniform chance that it will collide with a hydrogen nucleus. At first, we thought that strangeness conservation would prevent the K-particle from making a Λ^0 in such an interaction. Now, however, we see that that is not right. For although our K-particle *starts out* as a K^0 —which cannot make a Λ^0 —it does not *stay* this way. After a while, there *is some amplitude* that it will have flipped to the \bar{K}^0 state. We can, therefore, sometimes expect to see a Λ^0 produced along the K-particle track. The chance of this happening is given by the amplitude C_- , which we can [by using Eq. (11.50) backwards] relate to C_1 and C_2 . The relation is

$$C_- = \frac{1}{\sqrt{2}} (C_1 - C_2) = \frac{1}{2}(e^{-\beta t} e^{-i\alpha t} - 1). \quad (11.55)$$

As our K-particle goes along, the probability that it will “act like” a \bar{K}^0 is equal to $|C_-|^2$, which is

$$|C_-|^2 = \frac{1}{4}(1 + e^{-2\beta t} - 2e^{-\beta t} \cos \alpha t). \quad (11.56)$$

A complicated and strange result!

This, then, is the remarkable prediction of Gell-Mann and Pais: when a K^0 is produced, the chance that it will turn into a \bar{K}^0 —as it can demonstrate by being able to produce a Λ^0 —varies with time according to Eq. (11.56). This prediction came from using only sheer logic and the basic principles of the quantum mechanics—with no knowledge at all of the inner workings of the K-particle. Since nobody knows anything about the inner machinery, that is as far as Gell-Mann and Pais could go. They could not give any theoretical values for α and β . And nobody has been able to do so to this date. They were able to give a value of β obtained from the experimentally observed rate of decay into two π ’s ($2\beta = 10^{10}$ sec), but they could say nothing about α .

We have plotted the function of Eq. (11.56) for two values of α in Fig. 11-6. You can see that the form depends very much on the ratio of α to β . There is no \bar{K}^0 probability at first; then it builds up. If α is large, the probability would have

large oscillations. If α is small, there will be little or no oscillation—the probability will just rise smoothly to $1/4$.

Now, typically, the K-particle will be travelling at a constant speed near the speed of light. The curves of Fig. 11-6 then also represent the probability along the track of observing a \bar{K}^0 —with typical distances of several centimeters. You can see why this prediction is so remarkably peculiar. You produce a single particle and instead of just disintegrating, it does something else. Sometimes it disintegrates, and other times it turns into a different kind of a particle. Its characteristic probability of producing an effect varies in a strange way as it goes along. There is nothing else quite like it in nature. And this most remarkable prediction was made solely by arguments about the interference of amplitudes.

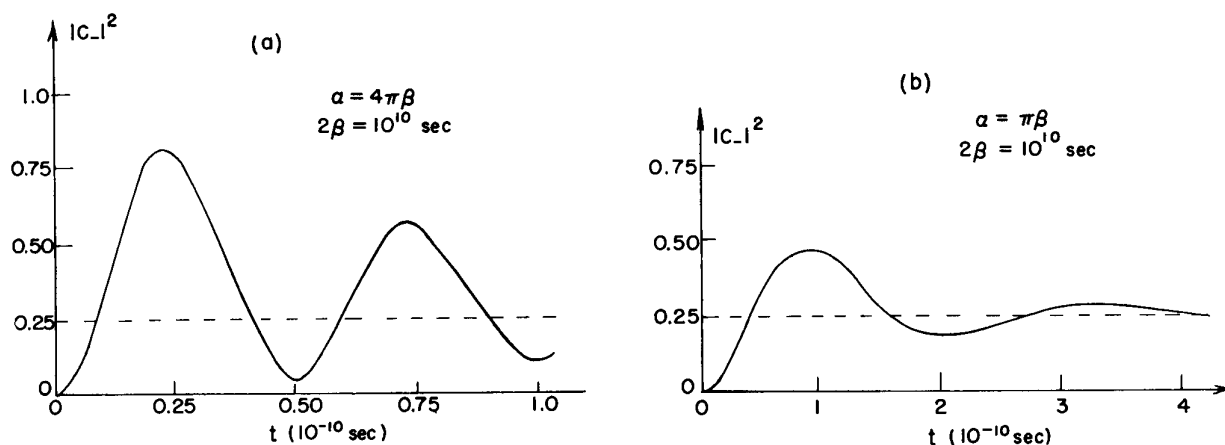


Fig. 11-6. The function of Eq. (11-56): (a) for $\alpha = \pi\beta$, (b) for $\alpha = 4\pi\beta$ (with $2\beta = 10^{10}$ sec).

If there is any place where we have a chance to test the main principles of quantum mechanics in the purest way—does the superposition of amplitudes work or doesn't it?—this is it. In spite of the fact that this effect has been predicted now for several years, there is no experimental determination that is very clear. There are some rough results which indicate that the α is not zero, and that the effect really occurs—they indicate that α is between 2β and 4β . That's all there is, experimentally. It would be very beautiful to check out the curve exactly to see if the principle of superposition really still works in such a mysterious world as that of the strange particles—with unknown reasons for the decays, and unknown reasons for the strangeness.

The analysis we have just described is very characteristic of the way quantum mechanics is being used today in the search for an understanding of the strange particles. All the complicated theories that you may hear about are no more and no less than this kind of elementary hocus-pocus using the principles of superposition and other principles of quantum mechanics of that level. Some people claim that they have theories by which it is possible to calculate the β and α , or at least the α given the β , but these theories are completely useless. For instance, the theory that predicts the value of α , given the β , tells us that the value of α should be infinite. The set of equations with which they originally start involves two π -mesons and then goes from the two π 's back to a K^0 , and so on. When it's all worked out, it does indeed produce a pair of equations like the ones we have here; but because there are an infinite number of states of two π 's, depending on their momenta, integrating over all the possibilities gives an α which is infinite. But nature's α is *not* infinite. So the dynamical theories are wrong. It is really quite remarkable that the phenomena which can be predicted *at all* in the world of the strange particles come from the principles of quantum mechanics at the level at which you are learning them now.

11-6 Generalization to N -state systems

We have finished with all the two-state systems we wanted to talk about. In the following chapters we will go on to study systems with more states. The extension to N -state systems of the ideas we have worked out for two states is pretty straightforward. It goes like this.

If a system has N distinct states, we can represent any state $|\psi(t)\rangle$ as a linear combination of any set of base states $|i\rangle$, where $i = 1, 2, 3, \dots, N$;

$$|\psi(t)\rangle = \sum_{\text{all } i} |i\rangle C_i(t). \quad (11.57)$$

The coefficients $C_i(t)$ are the amplitudes $\langle i | \psi(t) \rangle$. The behavior of the amplitudes C_i with time is governed by the equations

$$i\hbar \frac{dC_i(t)}{dt} = \sum_j H_{ij} C_j, \quad (11.58)$$

where the energy matrix H_{ij} describes the physics of the problem. It looks the same as for two states. Only now, both i and j must range over all N base states, and the energy matrix H_{ij} —or, if you prefer, the Hamiltonian—is an N by N matrix with N^2 numbers. As before, $H_{ij}^* = H_{ji}$ —so long as particles are conserved—and the diagonal elements H_{ii} are real numbers.

We have found a general solution for the C 's of a two-state system when the energy matrix is constant (doesn't depend on t). It is also not difficult to solve Eq. (11.58) for an N -state system when H is not time dependent. Again, we begin by looking for a possible solution in which the amplitudes all have the *same* time dependence. We try

$$C_i = a_i e^{-(i/\hbar)Et}. \quad (11.59)$$

When these C_i 's are substituted into (11.58), the derivatives $dC_i(t)/dt$ become just $(-i/\hbar)EC_i$. Canceling the common exponential factor from all terms, we get

$$Ea_i = \sum_j H_{ij} a_j. \quad (11.60)$$

This is a set of N linear algebraic equations for the N unknowns a_1, a_2, \dots, a_n , and there is a solution only if you are lucky—only if the determinant of the coefficients of all the a 's is zero. But it's not necessary to be that sophisticated; you can just start to solve the equations any way you want, and you will find that they can be solved only for certain values of E . (Remember that E is the only adjustable thing we have in the equations.)

If you want to be formal, however, you can write Eq. (11.60) as

$$\sum_j (H_{ij} - \delta_{ij}E) a_j = 0. \quad (11.61)$$

Then you can use the rule—if you know it—that these equations will have a solution only for those values of E for which

$$\text{Det } (H_{ij} - \delta_{ij}E) = 0. \quad (11.62)$$

Each term of the determinant is just H_{ij} , except that E is subtracted from every diagonal element. That is, (11.62) means just

$$\text{Det} \begin{pmatrix} H_{11} - E & H_{12} & H_{13} & \dots \\ H_{21} & H_{22} - E & H_{23} & \dots \\ H_{31} & H_{32} & H_{33} - E & \dots \\ \dots & \dots & \dots & \dots \end{pmatrix} = 0. \quad (11.63)$$

This is, of course, just a special way of writing an algebraic equation for E which is the sum of a bunch of products of all the terms taken a certain way. These products will give all the powers of E up to E^N .

So we have an N th order polynomial equal to zero, and there are, in general, N roots. (We must remember, however, that some of them may be multiple roots—meaning that two or more roots are equal.) Let's call the N roots

$$E_I, E_{II}, E_{III}, \dots, E_n, \dots, E_N. \quad (11.64)$$

(We will use n to represent the n th Roman numeral, so that n takes on the values I, II, \dots , N.) It may be that some of these energies are equal—say $E_{II} = E_{III}$ —but we will still choose to call them by different names.

The equations (11.60)—or (11.61)—have one solution for each value of E . If you put any one of the E 's—say E_n —into (11.60) and solve for the a_i , you get a set which belongs to the energy E_n . We will call this set $a_i(n)$.

Using these $a_i(n)$ in Eq. (11.59), we have the amplitudes $C_i(n)$ that the definite energy states are in the base state $|i\rangle$. Letting $|n\rangle$ stand for the state vector of the definite energy state at $t = 0$, we can write

$$C_i(n) = \langle i | n \rangle e^{(i/\hbar) E_n t},$$

with

$$\langle i | n \rangle = a_i(n). \quad (11.65)$$

The complete definite energy state $|\psi_n(t)\rangle$ can then be written as

$$|\psi_n(t)\rangle = \sum_i |i\rangle a_i(n) e^{-(i/\hbar) E_n t},$$

or

$$|\psi_n(t)\rangle = |n\rangle e^{(i/\hbar) E_n t}. \quad (11.66)$$

The state vectors $|n\rangle$ describe the configuration of the definite energy states, but have the time dependence factored out. Then they are constant vectors which can be used as a new base set if we wish.

Each of the states $|n\rangle$ has the property—as you can easily show—that when operated on by the Hamiltonian operator \hat{H} it gives just E_n times the same state:

$$\hat{H} |n\rangle = E_n |n\rangle. \quad (11.67)$$

The energy E_n is, then, a number which is a characteristic of the Hamiltonian operator \hat{H} . As we have seen, a Hamiltonian will, in general, have several characteristic energies. In the mathematician's world these would be called the “characteristic values” of the matrix H_{ij} . Physicists usually call them the “eigenvalues” of \hat{H} . (“Eigen” is the German word for “characteristic” or “proper.”) With each eigenvalue of \hat{H} —in other words, for each energy—there is the state of definite energy, which we have called the “stationary state.” Physicists usually call the states $|n\rangle$ “the eigenstates of \hat{H} .” Each eigenstate corresponds to a particular eigenvalue E_n .

Now, generally, the states $|n\rangle$ —of which there are N —can also be used as a base set. For this to be true, all of the states must be orthogonal, meaning that for any two of them, say $|n\rangle$ and $|m\rangle$,

$$\langle n | m \rangle = 0. \quad (11.68)$$

This will be true automatically if all the energies are different. Also, we can multiply all the $a_i(n)$ by a suitable factor so that all the states are normalized—by which we mean that

$$\langle n | n \rangle = 1 \quad (11.69)$$

for all n .

When it happens that Eq. (11.63) accidentally has two (or more) roots with the same energy, there are some minor complications. First, there are still two different sets of a_i 's which go with the two equal energies, but the states they give

may *not* be orthogonal. Suppose you go through the normal procedure and find two stationary states with equal energies—let's call them $|\mu\rangle$ and $|\nu\rangle$. Then it will not necessarily be so that they are orthogonal—if you are unlucky,

$$\langle\mu|\nu\rangle \neq 0.$$

It is, however, always true that you can cook up two new states, which we will call $|\mu'\rangle$ and $|\nu'\rangle$, that have the same energies and are also orthogonal, so that

$$\langle\mu'|\nu'\rangle = 0. \quad (11.70)$$

You can do this by making $|\mu'\rangle$ and $|\nu'\rangle$ a suitable linear combination of $|\mu\rangle$ and $|\nu\rangle$, with the coefficients chosen to make it come out so that Eq. (11.70) is true. It is always convenient to do this. We will generally assume that this has been done so that we can always assume that our proper energy states $|\mathbf{n}\rangle$ are all orthogonal.

We would like, for fun, to prove that when two of the stationary states have different energies they are indeed orthogonal. For the state $|\mathbf{n}\rangle$ with the energy $E_{\mathbf{n}}$, we have that

$$\hat{H}|\mathbf{n}\rangle = E_{\mathbf{n}}|\mathbf{n}\rangle. \quad (11.71)$$

This operator equation really means that there is an equation between numbers. Filling the missing parts, it means the same as

$$\sum_j \langle i|\hat{H}|j\rangle\langle j|\mathbf{n}\rangle = E_{\mathbf{n}}\langle i|\mathbf{n}\rangle. \quad (11.72)$$

If we take the complex conjugate of this equation, we get

$$\sum_j \langle i|\hat{H}|j\rangle^*\langle j|\mathbf{n}\rangle^* = E_{\mathbf{n}}^*\langle i|\mathbf{n}\rangle^*. \quad (11.73)$$

Remember now that the complex conjugate of an amplitude is the reverse amplitude, so (11.73) can be rewritten as

$$\sum_j \langle \mathbf{n}|j\rangle\langle j|\hat{H}|i\rangle = E_{\mathbf{n}}^*\langle \mathbf{n}|i\rangle. \quad (11.74)$$

Since this equation is valid for *any* i , its “short form” is

$$\langle \mathbf{n}|\hat{H} = E_{\mathbf{n}}^*\langle \mathbf{n}|, \quad (11.75)$$

which is called the *adjoint* to Eq. (11.71).

Now we can easily prove that $E_{\mathbf{n}}$ is a real number. We multiply Eq. (11.71) by $\langle \mathbf{n}|$ to get

$$\langle \mathbf{n}|\hat{H}|\mathbf{n}\rangle = E_{\mathbf{n}}, \quad (11.76)$$

since $\langle \mathbf{n}|\mathbf{n}\rangle = 1$. Then we multiply Eq. (11.75) on the left by $|\mathbf{n}\rangle$ to get

$$\langle \mathbf{n}|\hat{H}|\mathbf{n}\rangle = E_{\mathbf{n}}^*. \quad (11.77)$$

Comparing (11.76) with (11.77) it is clear that

$$E_{\mathbf{n}} = E_{\mathbf{n}}^*, \quad (11.78)$$

which means that $E_{\mathbf{n}}$ is real. We can erase the star on $E_{\mathbf{n}}$ in Eq. (11.75).

Finally we are ready to show that the different energy states are orthogonal. Let $|\mathbf{n}\rangle$ and $|\mathbf{m}\rangle$ be any two of the definite energy base states. Using Eq. (11.75) for the state \mathbf{m} , and multiplying it by $|\mathbf{n}\rangle$, we get that

$$\langle \mathbf{m}|\hat{H}|\mathbf{n}\rangle = E_{\mathbf{m}}\langle \mathbf{m}|\mathbf{n}\rangle.$$

But if we multiply (11.71) by $\langle \mathbf{m} |$, we get

$$\langle \mathbf{m} | \hat{H} | \mathbf{n} \rangle = E_{\mathbf{n}} \langle \mathbf{m} | \mathbf{n} \rangle.$$

Since the left sides of these two equations are equal, the right sides are, also:

$$E_{\mathbf{m}} \langle \mathbf{m} | \mathbf{n} \rangle = E_{\mathbf{n}} \langle \mathbf{m} | \mathbf{n} \rangle. \quad (11.79)$$

If $E_{\mathbf{m}} = E_{\mathbf{n}}$ the equation does not tell us anything. But if the energies of the two states $|\mathbf{m}\rangle$ and $|\mathbf{n}\rangle$ are *different* ($E_{\mathbf{m}} \neq E_{\mathbf{n}}$), Eq. (11.79) says that $\langle \mathbf{m} | \mathbf{n} \rangle$ must be zero, as we wanted to prove. The two states are necessarily orthogonal so long as $E_{\mathbf{n}}$ and $E_{\mathbf{m}}$ are numerically different.

The Hyperfine Splitting in Hydrogen

12-1 Base states for a system with two spin one-half particles

In this chapter we take up the “hyperfine splitting” of hydrogen, because it is a physically interesting example of what we can already do with quantum mechanics. It’s an example with more than two states, and it will be illustrative of the methods of quantum mechanics as applied to slightly more complicated problems. It is enough more complicated that once you see how this one is handled you can get immediately the generalization to all kinds of problems.

As you know, the hydrogen atom consists of an electron sitting in the neighborhood of the proton, where it can exist in any one of a number of discrete energy states in each one of which the pattern of motion of the electron is different. The first excited state, for example, lies $3/4$ of a Rydberg, or about 10 electron volts, above the ground state. But even the so-called ground state of hydrogen is not really a single, definite-energy state, because of the spins of the electron and the proton. These spins are responsible for the “hyperfine structure” in the energy levels, which splits all the energy levels into several nearly equal levels.

The electron can have its spin either “up” or “down” and, the proton can also have *its* spin either “up” or “down.” There are, therefore, *four* possible spin states for every dynamical condition of the atom. That is, when people say “the ground state” of hydrogen, they really mean the “four ground states,” and not just the very lowest state. The four spin states do not all have exactly the same energy; there are slight shifts from the energies we would expect with no spins. The shifts are, however, much, much smaller than the 10 volts or so from the ground state to the next state above. As a consequence, each dynamical state has its energy split into a set of very close energy levels—the so-called *hyperfine splitting*.

The energy differences among the four spin states is what we want to calculate in this chapter. The hyperfine splitting is due to the interaction of the magnetic moments of the electron and proton, which gives a slightly different magnetic energy for each spin state. These energy shifts are only about ten-millionths of an electron volt—really very small compared with 10 volts! It is because of this large gap that we can think about the ground state of hydrogen as a “four-state” system, without worrying about the fact that there are really many more states at higher energies. We are going to limit ourselves here to a study of the hyperfine structure of the ground state of the hydrogen atom.

For our purposes we are not interested in any of the details about the *positions* of the electron and proton because that has all been worked out by the atom so to speak—it has worked itself out by getting into the ground state. We need know only that we have an electron and proton in the neighborhood of each other with some definite spatial relationship. In addition, they can have various different relative orientations of their spins. It is only the effect of the spins that we want to look into.

The first question we have to answer is: What are the *base states* for the system? Now the question has been put incorrectly. There is no such thing as “*the*” base states, because, of course, the set of base states you may choose is not unique. New sets can always be made out of linear combinations of the old. There are always many choices for the base states, and among them, any choice is equally legitimate. So the question is not what is *the* base set, but what *could* a base set be? We can choose any one we wish for our own convenience. It is usually best to start with a base set which is *physically* the clearest. It may not be the solution

12-1 Base states for a system with two spin one-half particles

12-2 The Hamiltonian for the ground state of hydrogen

12-3 The energy levels

12-4 The Zeeman splitting

12-5 The states in a magnetic field

12-6 The projection matrix for spin one

to any problem, or may not have any *direct* importance, but it will generally make it easier to understand what is going on.

We choose the following four base states:

State 1: The electron and proton are both spin “up.”

State 2: The electron is “up” and the proton is “down.”

State 3: The electron is “down” and the proton is “up.”

State 4: The electron and proton are both “down.”

We need a handy notation for these four states, so we’ll represent them this way:

State 1: $| + + \rangle$; electron *up*, proton *up*.

State 2: $| + - \rangle$; electron *up*, proton *down*.

State 3: $| - + \rangle$; electron *down*, proton *up*.

State 4: $| - - \rangle$; electron *down*, proton *down*.

(12.1)

You will have to remember that the *first* plus or minus sign refers to the electron and the *second*, to the proton. For handy reference, we’ve also summarized the notation in Fig. 12-1. Sometimes it will also be convenient to call these states $| 1 \rangle$, $| 2 \rangle$, $| 3 \rangle$, and $| 4 \rangle$.

You may say, “But the particles interact, and maybe these aren’t the right base states. It sounds as though you are considering the two particles independently.” Yes, indeed! The interaction raises the problem: what is the *Hamiltonian* for the system, but the interaction is not involved in the question of how to *describe* the system. What we choose for the base states has nothing to do with what happens next. It may be that the atom cannot ever *stay* in one of these base states, even if it is started that way. That’s another question. That’s the question: How do the amplitudes change with time in a particular (fixed) base? In choosing the base states, we are just choosing the “unit vectors” for our description.

While we’re on the subject, let’s look at the general problem of finding a set of base states when there is more than one particle. You know the base states for a single particle. An electron, for example, is completely described in real life—not in our simplified cases, but in real life—by giving the amplitudes to be in each of the following states:

$| \text{electron “up” with momentum } \mathbf{p} \rangle$

or

$| \text{electron “down” with momentum } \mathbf{p} \rangle$.

There are really two infinite sets of states, one state for each value of \mathbf{p} . That is to say that an electron state $|\psi\rangle$ is completely described if you know all the amplitudes

$$\langle +, \mathbf{p} | \psi \rangle \quad \text{and} \quad \langle -, \mathbf{p} | \psi \rangle,$$

where the $+$ and $-$ represent the components of angular momentum along some axis—usually the z -axis—and \mathbf{p} is the vector momentum. There must, therefore, be two amplitudes for every possible momentum (a multi-infinite set of base states). That is all there is to describing a single particle.

When there is more than one particle, the base states can be written in a similar way. For instance, if there were an electron and a proton in a more complicated situation than we are considering, the base states could be of the following kind:

$| \text{an electron with spin “up,” moving with momentum } \mathbf{p}_1 \text{ and}$

$\text{a proton with spin “down,” moving with momentum } \mathbf{p}_2 \rangle$.

And so on for other spin combinations. If there are more than two particles—same idea. So you see that to write down the *possible* base states is really very easy. The only problem is, what is the Hamiltonian?

For our study of the ground state of hydrogen we don’t need to use the full sets of base states for the various momenta. We are specifying particular mo-

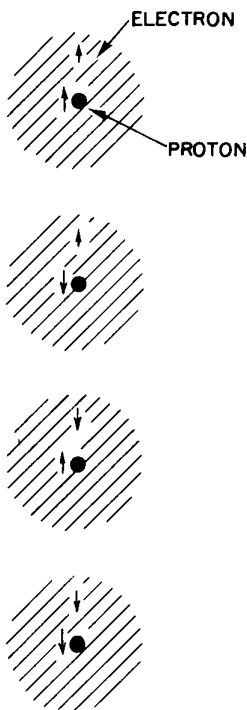


Fig. 12-1. A set of base states for the ground state of the hydrogen atom.

mentum states for the proton and electron when we say “the ground state.” The details of the configuration—the amplitudes for all the momentum base states—can be calculated, but that is another problem. Now we are concerned only with the effects of the spin, so we can take only the four base states of (12.1). Our next problem is: What is the Hamiltonian for this set of states?

12-2 The Hamiltonian for the ground state of hydrogen

We’ll tell you in a moment what it is. But first, we should remind you of one thing: *any* state can always be written as a linear combination of the base states. For any state $|\psi\rangle$ we can write

$$|\psi\rangle = |++\rangle\langle++|\psi\rangle + |+-\rangle\langle+-|\psi\rangle + |-+\rangle\langle-+|\psi\rangle + |--\rangle\langle--|\psi\rangle. \quad (12.2)$$

Remember that the complete brackets are just complex numbers, so we can also write them in the usual fashion as C_i , where $i = 1, 2, 3$, or 4 , and write Eq. (12.2) as

$$|\psi\rangle = |++\rangle C_1 + |+-\rangle C_2 + |-+\rangle C_3 + |--\rangle C_4. \quad (12.3)$$

By giving the four amplitudes C_i we completely describe the spin state $|\psi\rangle$. If these four amplitudes change with time, as they will, the rate of change in time is given by the operator \hat{H} . The problem is to find the \hat{H} .

There is no general rule for writing down the Hamiltonian of an atomic system, and finding the right formula is much more of an art than finding a set of base states. We were able to tell you a general rule for writing a set of base states for any problem of a proton and an electron, but to describe the general Hamiltonian of such a combination is too hard at this level. Instead, we will lead you to a Hamiltonian by some heuristic argument—and you will have to accept it as the correct one because the results will agree with the test of experimental observation.

You will remember that in the last chapter we were able to describe the Hamiltonian of a single, spin one-half particle by using the sigma matrices—or the exactly equivalent sigma operators. The properties of the operators are summarized in Table 12-1. These operators—which are just a convenient, shorthand way of keeping track of the matrix elements of the type $\langle+|\sigma_z|+\rangle$ —were useful for describing the behavior of a *single* particle of spin one-half. The question is: Can we find an analogous device to describe a system with two spins? The answer is yes, very simply, as follows. We invent a thing which we will call “sigma electron,” which we represent by the vector operator σ^e , and which has the x -, y -, and z -components, $\sigma_x^e, \sigma_y^e, \sigma_z^e$. We now make the *convention* that when one of these things operates on any one of our four base states of the hydrogen atom, it acts only on the *electron* spin, and in exactly the same way as if the electron were all by itself. Example: What is $\sigma_y^e| - + \rangle$? Since σ_y on an electron “down” is $-i$ times the corresponding state with the electron “up”,

$$\sigma_y^e| - + \rangle = -i| + + \rangle.$$

(When σ_y^e acts on the combined state it flips over the electron, but does nothing to the proton and multiplies the result by $-i$.) Operating on the other states, σ_y^e would give

$$\begin{aligned} \sigma_y^e| + + \rangle &= i| - + \rangle, \\ \sigma_y^e| + - \rangle &= i| - - \rangle, \\ \sigma_y^e| - - \rangle &= -i| + - \rangle. \end{aligned}$$

Just remember that the operators σ^e work only on the *first* spin symbol—that is, on the *electron* spin.

Next we define the corresponding operator “sigma proton” for the proton spin. Its three components $\sigma_x^p, \sigma_y^p, \sigma_z^p$ act in the same way as σ^e , only on the

Table 12-1

$\sigma_z +\rangle$	$= + +\rangle$
$\sigma_z -\rangle$	$= - -\rangle$
$\sigma_x +\rangle$	$= + -\rangle$
$\sigma_x -\rangle$	$= + +\rangle$
$\sigma_y +\rangle$	$= +i -\rangle$
$\sigma_y -\rangle$	$= -i +\rangle$

proton spin. For example, if we have σ_x^p acting on each of the four base states, we get—always using Table 12-1—

$$\begin{aligned}\sigma_x^p | + + \rangle &= | + - \rangle, \\ \sigma_x^p | + - \rangle &= | + + \rangle, \\ \sigma_x^p | - + \rangle &= | - - \rangle, \\ \sigma_x^p | - - \rangle &= | - + \rangle.\end{aligned}$$

As you can see, it's not very hard.

Now in the most general case we could have more complex things. For instance, we could have products of the two operators like $\sigma_y^e \sigma_z^p$. When we have such a product we do first what the operator on the right says, and then do what the other one says.† For example, we would have that

$$\sigma_x^e \sigma_z^p | + - \rangle = \sigma_x^e (\sigma_z^p | + - \rangle) = \sigma_x^e (- | + - \rangle) = -\sigma_x^e | + - \rangle = - | - - \rangle.$$

Note that these operators don't do anything on pure numbers—we have used this fact when we wrote $\sigma_x^e(-1) = (-1)\sigma_x^e$. We say that the operators “commute” with pure numbers, or that a number “can be moved through” the operator. You can practice by showing that the product $\sigma_x^e \sigma_z^p$ gives the following results for the four states:

$$\begin{aligned}\sigma_x^e \sigma_z^p | + + \rangle &= + | - + \rangle, \\ \sigma_x^e \sigma_z^p | + - \rangle &= - | - - \rangle, \\ \sigma_x^e \sigma_z^p | - + \rangle &= + | + + \rangle, \\ \sigma_x^e \sigma_z^p | - - \rangle &= - | + - \rangle.\end{aligned}$$

If we take all the possible operators, using each kind of operator only once, there are sixteen possibilities. Yes, *sixteen*—provided we include also the “unit operator” $\hat{1}$. First, there are the three: $\sigma_x^e, \sigma_y^e, \sigma_z^e$. Then the three $\sigma_z^p, \sigma_y^p, \sigma_x^p$ —that makes six. In addition, there are the nine possible products of the form $\sigma_x^e \sigma_y^p$, which makes a total of 15. And there's the unit operator which just leaves any state unchanged. Sixteen in all.

Now note that for a four-state system, the Hamiltonian matrix has to be a four-by-four matrix of coefficients—it will have sixteen entries. It is easily demonstrated that any four-by-four matrix—and, therefore, the Hamiltonian matrix in particular—can be written as a linear combination of the sixteen double-spin matrices corresponding to the set of operators we have just made up. Therefore, for the interaction between a proton and an electron that involves only their spins, we can expect that the Hamiltonian operator can be written as a linear combination of the same 16 operators. The only question is, how?

Well, first, we know that the interaction doesn't depend on our choice of axes for a coordinate system. If there is no external disturbance—like a magnetic field—to determine a unique direction in space, the Hamiltonian can't depend on our choice of the direction of the x -, y -, and z -axes. That means that the Hamiltonian can't have a term like σ_x^e all by itself. It would be ridiculous, because then somebody with a different coordinate system would get different results.

The only possibilities are a term with the unit matrix, say a constant a (times $\hat{1}$), and some combination of the sigmas that doesn't depend on the coordinates—some “invariant” combination. The only *scalar* invariant combination of two vectors is the dot product, which for our σ 's is

$$\sigma^e \cdot \sigma^p = \sigma_x^e \sigma_x^p + \sigma_y^e \sigma_y^p + \sigma_z^e \sigma_z^p. \quad (12.4)$$

This operator is invariant with respect to any rotation of the coordinate system.

† For these *particular* operators, you will notice it turns out that the sequence of the operators doesn't matter.

So the only possibility for a Hamiltonian with the proper symmetry in space is a constant times the unit matrix plus a constant times this dot product, say,

$$\hat{H} = E_0 + A \boldsymbol{\sigma}^e \cdot \boldsymbol{\sigma}^p. \quad (12.5)$$

That's our Hamiltonian. It's the only thing that it can be, by the symmetry of space, *so long as there is no external field*. The constant term doesn't tell us much; it just depends on the level we choose to measure energies from. We may just as well take $E_0 = 0$. The second term tells us all we need to know to find the level splitting of the hydrogen.

If you want to, you can think of the Hamiltonian in a different way. If there are two magnets near each other with magnetic moments μ_e and μ_p , the mutual energy will depend on $\mu_e \cdot \mu_p$ —among other things. And, you remember, we found that the classical thing we call μ_e appears in quantum mechanics as $\mu_e \sigma_e$. Similarly, what appears classically as μ_p will usually turn out in quantum mechanics to be $\mu_p \sigma_p$ (where μ_p is the magnetic moment of the proton, which is about 1000 times smaller than μ_e , and has the opposite sign). So Eq. (12.5) says that the interaction energy is like the interaction between two magnets—only not quite, because the interaction of the two magnets depends on the radial distance between them. But Eq. (12.5) could be—and, in fact, *is*—some kind of an average interaction. The electron is moving all around inside the atom, and our Hamiltonian gives only the average interaction energy. All it says is that for a prescribed arrangement in space for the electron and proton there is an energy proportional to the cosine of the angle between the two magnetic moments, speaking classically. Such a classical qualitative picture may help you to understand where it comes from, but the important thing is that Eq. (12.5) is the correct quantum mechanical formula.

The order of magnitude of the classical interaction between two magnets would be the product of the two magnetic moments divided by the cube of the distance between them. The distance between the electron and the proton in the hydrogen atom is, speaking roughly, one half an atomic radius, or 0.5 angstrom. It is, therefore, possible to make a crude estimate that the constant A should be about equal to the product of the two magnetic moments μ_e and μ_p divided by the cube of 1/2 angstrom. Such an estimate gives a number in the right ball park. It turns out that A can be calculated accurately once you understand the complete quantum theory of the hydrogen atom—which we so far do not. It has, in fact, been calculated to an accuracy of about 30 parts in one million. So, unlike the flip-flop constant A of the ammonia molecule, which couldn't be calculated at all well by a theory, our constant A for the hydrogen *can* be calculated from a more detailed theory. But never mind, we will for our present purposes think of the A as a number which could be determined by experiment, and analyze the physics of the situation.

Taking the Hamiltonian of Eq. (12.5), we can use it with the equation

$$i\hbar \dot{C}_i = \sum_j H_{ij} C_j \quad (12.6)$$

to find out what the spin interactions do to the energy levels. To do that, we need to work out the sixteen matrix elements $H_{ij} = \langle i | H | j \rangle$ corresponding to each pair of the four base states in (12.1).

We begin by working out what $\hat{H} |j\rangle$ is for each of the four base states. For example,

$$\hat{H} |++\rangle = A \boldsymbol{\sigma}^e \cdot \boldsymbol{\sigma}^p |++\rangle = A \{ \sigma_x^e \sigma_x^p + \sigma_y^e \sigma_y^p + \sigma_z^e \sigma_z^p \} |++\rangle. \quad (12.7)$$

Using the method we described a little earlier—it's easy if you have memorized Table 12-1—we find what each pair of σ 's does on $|++\rangle$. The answer is

$$\begin{aligned} \sigma_x^e \sigma_x^p |++\rangle &= + |--\rangle, \\ \sigma_y^e \sigma_y^p |++\rangle &= - |--\rangle, \\ \sigma_z^e \sigma_z^p |++\rangle &= + |++\rangle. \end{aligned} \quad (12.8)$$

So (12.7) becomes

$$\hat{H} | + + \rangle = A \{ | - - \rangle - | - + \rangle + | + + \rangle \} = A | + + \rangle. \quad (12.9)$$

Since our four base states are all orthogonal, that gives us immediately that

$$\begin{aligned} \langle + + | \hat{H} | + + \rangle &= A \langle + + | + + \rangle = A, \\ \langle + - | \hat{H} | + + \rangle &= A \langle + - | + + \rangle = 0, \\ \langle - + | \hat{H} | + + \rangle &= A \langle - + | + + \rangle = 0, \\ \langle - - | \hat{H} | + + \rangle &= A \langle - - | + + \rangle = 0. \end{aligned} \quad (12.10)$$

Remembering that $\langle j | \hat{H} | i \rangle = \langle i | \hat{H} | j \rangle^*$, we can already write down the differential equation for the amplitudes C_1 :

$$i\hbar \dot{C}_1 = H_{11}C_1 + H_{12}C_2 + H_{13}C_3 + H_{14}C_4$$

or

$$i\hbar \dot{C}_1 = AC_1. \quad (12.11)$$

That's all! We get only the one term.

Now to get the rest of the Hamiltonian equations we have to crank through the same procedure for \hat{H} operating on the other states. First, we will let you practice by checking out all of the sigma products we have written down in Table 12-2. Then we can use them to get:

$$\begin{aligned} \hat{H} | + - \rangle &= A \{ 2 | - + \rangle - | + - \rangle \}, \\ \hat{H} | - + \rangle &= A \{ 2 | + - \rangle - | - + \rangle \}, \\ \hat{H} | - - \rangle &= A | - - \rangle. \end{aligned} \quad (12.12)$$

Then, multiplying each one in turn on the left by all the other state vectors, we get the following Hamiltonian matrix, H_{ij} :

$$H_{ij} = \begin{matrix} & \begin{matrix} \xrightarrow{j} \\ \downarrow i \end{matrix} & \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \end{matrix} & \begin{pmatrix} A & 0 & 0 & 0 \\ 0 & -A & 2A & 0 \\ 0 & 2A & -A & 0 \\ 0 & 0 & 0 & A \end{pmatrix} \end{matrix}. \quad (12.13)$$

It means, of course, nothing more than that our differential equations for the four amplitudes C_i are

$$\begin{aligned} i\hbar \dot{C}_1 &= AC_1, \\ i\hbar \dot{C}_2 &= -AC_2 + 2AC_3, \\ i\hbar \dot{C}_3 &= 2AC_2 - AC_3, \\ i\hbar \dot{C}_4 &= AC_4. \end{aligned} \quad (12.14)$$

Before solving these equations we can't resist telling you about a clever rule due to Dirac—it will make you feel that you are really advanced—although we don't need it for our work. We have—from the equations (12.9) and (12.12)—that

$$\begin{aligned} \sigma^e \cdot \sigma^p | + + \rangle &= | + + \rangle, \\ \sigma^e \cdot \sigma^p | + - \rangle &= 2 | - + \rangle - | + - \rangle, \\ \sigma^e \cdot \sigma^p | - + \rangle &= 2 | + - \rangle - | - + \rangle, \\ \sigma^e \cdot \sigma^p | - - \rangle &= | - - \rangle. \end{aligned} \quad (12.15)$$

Table 12-2

Spin operators for the hydrogen atom

$\sigma_x^e \sigma_x^p + + \rangle$	$= + - - \rangle$
$\sigma_x^e \sigma_x^p + - \rangle$	$= + - + \rangle$
$\sigma_x^e \sigma_x^p - + \rangle$	$= + + - \rangle$
$\sigma_x^e \sigma_x^p - - \rangle$	$= + + + \rangle$
$\sigma_y^e \sigma_y^p + + \rangle$	$= - - - \rangle$
$\sigma_y^e \sigma_y^p + - \rangle$	$= + - + \rangle$
$\sigma_y^e \sigma_y^p - + \rangle$	$= + + - \rangle$
$\sigma_y^e \sigma_y^p - - \rangle$	$= - + + \rangle$
$\sigma_z^e \sigma_z^p + + \rangle$	$= + + + \rangle$
$\sigma_z^e \sigma_z^p + - \rangle$	$= - + - \rangle$
$\sigma_z^e \sigma_z^p - + \rangle$	$= - - + \rangle$
$\sigma_z^e \sigma_z^p - - \rangle$	$= + - - \rangle$

Look, said Dirac, I can also write the first and last equations as

$$\begin{aligned}\sigma^e \cdot \sigma^p |++\rangle &= 2|++\rangle - |++\rangle, \\ \sigma^e \cdot \sigma^p |--\rangle &= 2|--\rangle - |--\rangle;\end{aligned}$$

then they are all quite similar. Now I invent a new operator, which I will call $P_{\text{spin exch}}$ and which I *define* to have the following properties:†

$$\begin{aligned}P_{\text{spin exch}} |++\rangle &= |--\rangle, \\ P_{\text{spin exch}} |+-\rangle &= |-+\rangle, \\ P_{\text{spin exch}} |-+\rangle &= |+-\rangle, \\ P_{\text{spin exch}} |--\rangle &= |++\rangle.\end{aligned}$$

All the operator does is interchange the spin directions of the two particles. Then I can write the whole set of equations in (12.15) as a simple operator equation:

$$\sigma^e \cdot \sigma^p = 2P_{\text{spin exch}} - 1. \quad (12.16)$$

That's the formula of Dirac. His "spin exchange operator" gives a handy rule for figuring out $\sigma^e \cdot \sigma^p$. (You see, you can do everything now. The gates are open.)

12-3 The energy levels

Now we are ready to work out the energy levels of the ground state of hydrogen by solving the Hamiltonian equations (12.14). We want to find the energies of the stationary states. This means that we want to find those special states $|\psi\rangle$ for which each amplitude $C_i = \langle i | \psi \rangle$ in the set belonging to $|\psi\rangle$ has the same time dependence—namely, $e^{-i\omega t}$. Then the state will have the energy $E = \hbar\omega$. So we want a set for which

$$C_i = a_i e^{(-i/\hbar)Et}, \quad (12.17)$$

where the four coefficients a_i are independent of time. To see whether we can get such amplitudes, we substitute (12.17) into Eq. (12.14) and see what happens. Each $i\hbar dC/dt$ in Eq. (12.14) turns into EC , and—after cancelling out the common exponential factor—each C becomes an a ; we get

$$\begin{aligned}Ea_1 &= Aa_1, \\ Ea_2 &= -Aa_2 + 2Aa_3, \\ Ea_3 &= 2Aa_2 - Aa_3, \\ Ea_4 &= Aa_4,\end{aligned} \quad (12.18)$$

which we have to solve for a_1, a_2, a_3 , and a_4 . Isn't it nice that the first equation is independent of the rest—that means we can see one solution right away. If we choose $E = A$,

$$a_1 = 1, \quad a_2 = a_3 = a_4 = 0,$$

gives a solution. (Of course, taking all the a 's equal to zero also gives a solution, but that's no state at all!) Let's call our first solution the state $|I\rangle$:‡

$$|I\rangle = |I\rangle = |++\rangle. \quad (12.19)$$

Its energy is

$$E_I = A.$$

† This operator is now called the "Pauli spin exchange operator."

‡ The state is really $|I\rangle e^{(-i/\hbar)E_I t}$; but, as usual we will identify the states by the constant vectors which are equal to the complete vectors at $t = 0$.

With that clue you can immediately see another solution from the last equation in (12.18):

$$a_1 = a_2 = a_3 = 0, \quad a_4 = 1, \\ E = A.$$

We'll call that solution state $|II\rangle$:

$$|II\rangle = |4\rangle = |--\rangle, \\ E_{II} = A. \quad (12.20)$$

Now it gets a little harder; the two equations left in (12.18) are mixed up. But we've done it all before. Adding the two, we get

$$E(a_2 + a_3) = A(a_2 + a_3). \quad (12.21)$$

Subtracting, we have

$$E(a_2 - a_3) = -3A(a_2 - a_3). \quad (12.22)$$

By inspection—and remembering ammonia—we see that there are two solutions:

$$a_2 = a_3, \quad E = A \\ \text{and} \quad (12.23)$$

$$a_2 = -a_3, \quad E = -3A.$$

They are mixtures of $|2\rangle$ and $|3\rangle$. Calling these states $|III\rangle$ and $|IV\rangle$, and putting in a factor $1/\sqrt{2}$ to make the states properly normalized, we have

$$|III\rangle = \frac{1}{\sqrt{2}}(|2\rangle + |3\rangle) = \frac{1}{\sqrt{2}}(|+-\rangle + |-+\rangle), \quad (12.24)$$

$$E_{III} = A$$

and

$$|IV\rangle = \frac{1}{\sqrt{2}}(|2\rangle - |3\rangle) = \frac{1}{\sqrt{2}}(|+-\rangle - |-+\rangle), \quad (12.25)$$

$$E_{IV} = -3A.$$

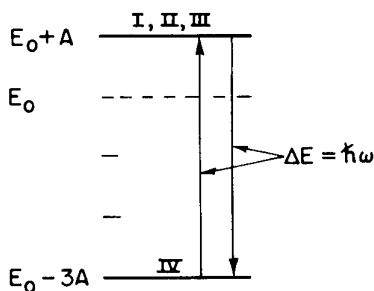


Fig. 12-2. Energy-level diagram for the ground state of atomic hydrogen.

We have found four stationary states and their energies. Notice, incidentally, that our four states are orthogonal, so they also can be used for base states if desired. Our problem is completely solved.

Three of the states have the energy A , and the last has the energy $-3A$. The average is zero—which means that when we took $E_0 = 0$ in Eq. (12.5), we were choosing to measure all the energies from the average energy. We can draw the energy-level diagram for the ground state of hydrogen as shown in Fig. 12-2.

Now the difference in energy between state $|IV\rangle$ and any one of the others is $4A$. An atom which happens to have gotten into state $|I\rangle$ could fall from there to state $|IV\rangle$ and emit light. Not optical light, because the energy is so tiny—it would emit a microwave quantum. Or, if we shine microwaves on hydrogen gas, we will find an absorption of energy as the atoms in state $|IV\rangle$ pick up energy and go into one of the upper states—but only at the frequency $\omega = 4A/\hbar$. This frequency has been measured experimentally; the best result, obtained very recently,[†] is

$$f = \omega/2\pi = (1,420,405,751.800 \pm 0.028) \text{ cycles per second.} \quad (12.26)$$

The error is only two parts in 100 billion! Probably no basic physical quantity is measured better than that—it's one of the most remarkably accurate measurements in physics. The theorists were very happy that they could compute the energy to an accuracy of 3 parts in 10^5 , but in the meantime it has been measured to 2 parts in 10^{11} —a million times more accurate than the theory. So the experimenters are

[†] Crampton, Kleppner, and Ramsey; *Physical Review Letters*, Vol. 11, page 338 (1963).

way ahead of the theorists. In the theory of the ground state of the hydrogen atom *you* are as good as anybody. You, too, can just take your value of A from experiment—that's what everybody has to do in the end.

You have probably heard before about the "21-centimeter line" of hydrogen. That's the wavelength of the 1420 megacycle spectral line between the hyperfine states. Radiation of this wavelength is emitted or absorbed by the atomic hydrogen gas in the galaxies. So with radio telescopes tuned in to 21-cm waves (or 1420 megacycles approximately) we can observe the velocities and the location of concentrations of atomic hydrogen gas. By measuring the intensity, we can estimate the amount of hydrogen. By measuring the frequency shift due to the Doppler effect, we can find out about the motion of the gas in the galaxy. That is one of the big programs of radio astronomy. So now we are talking about something that's very real—it is not an artificial problem.

12-4 The Zeeman splitting

Although we have finished the problem of finding the energy levels of the hydrogen ground state, we would like to study this interesting system some more. In order to say anything more about it—for instance, in order to calculate the rate at which the hydrogen atom absorbs or emits radio waves at 21 centimeters—we have to know what happens when the atom is disturbed. We have to do as we did for the ammonia molecule—after we found the energy levels we went on and studied what happened when the molecule was in an electric field. We were then able to figure out the effects from the electric field in a radio wave. For the hydrogen atom, the electric field does nothing to the levels, except to move them all by some constant amount proportional to the square of the field—which is not of any interest because that won't change the energy *differences*. It is now the *magnetic* field which is important. So the next step is to write the Hamiltonian for a more complicated situation in which the atom sits in an external magnetic field.

What, then, is the Hamiltonian? We'll just tell you the answer, because we can't give you any "proof" except to say that this is the way the atom works.

The Hamiltonian is

$$\hat{H} = A(\sigma^e \cdot \sigma^p) - \mu_e \sigma^e \cdot \mathbf{B} - \mu_p \sigma^p \cdot \mathbf{B}. \quad (12.27)$$

It now consists of three parts. The first term $A\sigma^e \cdot \sigma^p$ represents the magnetic interaction between the electron and the proton—it is the same one that would be there if there were no magnetic field. This is the term we have already had; and the influence of the magnetic field on the constant A is negligible. The effect of the external magnetic field shows up in the last two terms. The second term, $-\mu_e \sigma^e \cdot \mathbf{B}$, is the energy the electron would have in the magnetic field if it were there alone.† In the same way, the last term $-\mu_p \sigma^p \cdot \mathbf{B}$, would have been the energy of a proton alone. Classically, the energy of the two of them together would be the sum of the two, and that works also quantum mechanically. In a magnetic field, the energy of interaction due to the magnetic field is just the sum of the energy of interaction of the electron with the external field, and of the proton with the field—both expressed in terms of the sigma operators. In quantum mechanics these terms are not really the energies, but thinking of the classical formulas for the energy is a way of remembering the rules for writing down the Hamiltonian. Anyway, the correct Hamiltonian is Eq. (12.27).

Now we have to go back to the beginning and do the problem all over again. Much of the work is, however, done—we need only to add the effects of the new terms. Let's take a constant magnetic field \mathbf{B} in the z -direction. Then we have to

† Remember that classically $U = -\boldsymbol{\mu} \cdot \mathbf{B}$, so the energy is lowest when the moment is along the field. For positive particles, the magnetic moment is parallel to the spin and for negative particles it is opposite. So in Eq. (12.27), μ_p is a *positive* number, but μ_e is a *negative* number.

add to our Hamiltonian operator \hat{H} the two new pieces—which we can call \hat{H}' :

$$\hat{H}' = -(\mu_e \sigma_z^e + \mu_p \sigma_z^p)B.$$

Using Table 12-1, we get right away that

$$\begin{aligned}\hat{H}' |++\rangle &= -(\mu_e + \mu_p)B |++\rangle, \\ \hat{H}' |+-\rangle &= -(\mu_e - \mu_p)B |+-\rangle, \\ \hat{H}' |-+\rangle &= -(-\mu_e + \mu_p)B |-+\rangle, \\ \hat{H}' |--\rangle &= (\mu_e + \mu_p)B |--\rangle.\end{aligned}\tag{12.28}$$

How very convenient! The \hat{H}' operating on each state just gives a number times that state. The matrix $\langle i | \hat{H}' | j \rangle$ has, therefore, only *diagonal* elements—we can just add the coefficients in (12.28) to the corresponding diagonal terms of (12.13), and the Hamiltonian equations of (12.14) become

$$\begin{aligned}i\hbar dC_1/dt &= \{A - (\mu_e + \mu_p)B\}C_1, \\ i\hbar dC_2/dt &= -\{A + (\mu_e - \mu_p)B\}C_2 + 2AC_3, \\ i\hbar dC_3/dt &= 2AC_2 - \{A - (\mu_e - \mu_p)B\}C_3, \\ i\hbar dC_4/dt &= \{A + (\mu_e + \mu_p)B\}C_4.\end{aligned}\tag{12.29}$$

The form of the equations is not different—only the coefficients. So long as B doesn't vary with time, we can continue as we did before. Substituting $C_i = a_i e^{-(i/\hbar)Et}$, we get—as a modification of (12.18)—

$$\begin{aligned}Ea_1 &= A - (\mu_e + \mu_p)B a_1, \\ Ea_2 &= -\{A + (\mu_e - \mu_p)B\}a_2 + 2Aa_3, \\ Ea_3 &= 2Aa_2 - \{A - (\mu_e - \mu_p)B\}a_3, \\ Ea_4 &= \{A + (\mu_e + \mu_p)B\}a_4.\end{aligned}\tag{12.30}$$

Fortunately, the first and fourth equations are still independent of the rest, so the same technique works again.

One solution is the state $|I\rangle$ for which $a_1 = 1, a_2 = a_3 = a_4 = 0$, or

$$|I\rangle = |1\rangle = |++\rangle,\tag{12.31}$$

with

$$E_I = A - (\mu_e + \mu_p)B.$$

Another is

$$|II\rangle = |4\rangle = |--\rangle,\tag{12.32}$$

with

$$E_{II} = A + (\mu_e + \mu_p)B.$$

A little more work is involved for the remaining two equations, because the coefficients of a_2 and a_3 are no longer equal. But they are just like the pair we had for the ammonia molecule. Looking back at Eq. (9.20), we can make the following analogy (remembering that the labels 1 and 2 there correspond to 2 and 3 here):

$$\begin{aligned}H_{11} &\rightarrow -A - (\mu_e - \mu_p)B, \\ H_{12} &\rightarrow 2A, \\ H_{21} &\rightarrow 2A, \\ H_{22} &\rightarrow -A + (\mu_e - \mu_p)B.\end{aligned}\tag{12.33}$$

The energies are then given by (9.25), which was

$$E = \frac{H_{11} + H_{22}}{2} \pm \sqrt{\frac{(H_{11} - H_{22})^2}{4} + H_{12}H_{21}}.\tag{12.34}$$

Making the substitutions from (12.33), the energy formula becomes

$$E = -A \pm \sqrt{(\mu_e - \mu_p)^2 B^2 + 4A^2}.$$

Although in Chapter 9 we used to call these energies E_I and E_{II} , and we are in this problem calling them E_{III} and E_{IV} ,

$$\begin{aligned} E_{III} &= A\{-1 + 2\sqrt{1 + (\mu_e - \mu_p)^2 B^2 / 4A^2}\}, \\ E_{IV} &= -A\{1 + 2\sqrt{1 + (\mu_e - \mu_p)^2 B^2 / 4A^2}\}. \end{aligned} \quad (12.35)$$

So we have found the energies of the four stationary states of a hydrogen atom in a constant magnetic field. Let's check our results by letting B go to zero and seeing whether we get the same energies we had in the preceding section. You see that we do. For $B = 0$, the energies E_I , E_{II} , and E_{III} go to $+A$, and E_{IV} goes to $-3A$. Even our labeling of the states agrees with what we called them before. When we turn on the magnetic field though, all of the energies change in a different way. Let's see how they go.

First, we have to remember that for the electron, μ_e is negative, and about 1000 times larger than μ_p —which is positive. So $\mu_e + \mu_p$ and $\mu_e - \mu_p$ are both negative numbers, and nearly equal. Let's call them $-\mu$ and $-\mu'$:

$$\mu = -(\mu_e + \mu_p), \quad \mu' = -(\mu_e - \mu_p). \quad (12.36)$$

(Both μ and μ' are positive numbers, nearly equal to magnitude of μ_e —which is about one Bohr magneton.) Then our four energies are

$$\begin{aligned} E_I &= A + \mu B, \\ E_{II} &= A - \mu B, \\ E_{III} &= A\{-1 + 2\sqrt{1 + \mu'^2 B^2 / 4A^2}\}, \\ E_{IV} &= -A\{1 + 2\sqrt{1 + \mu'^2 B^2 / 4A^2}\}. \end{aligned} \quad (12.37)$$

The energy E_I starts at A and increases linearly with B —with the slope μ . The

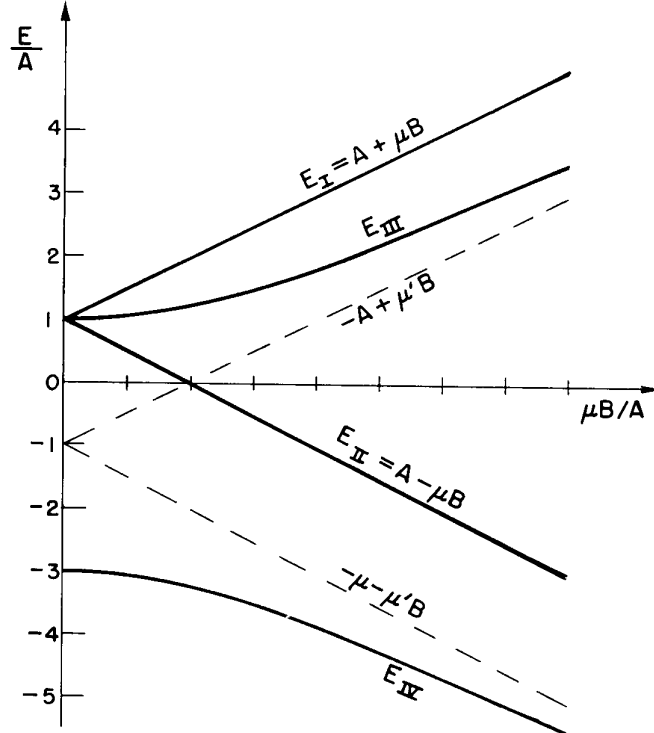


Fig. 12-3. The energy levels of the ground state of hydrogen in a magnetic field B .

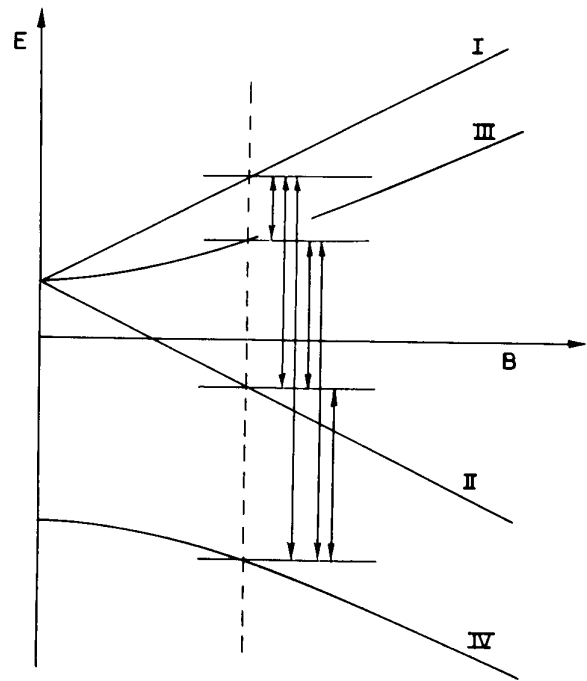


Fig. 12-4. Transitions between the levels of ground state energy levels of hydrogen in some particular magnetic field B .

energy E_{II} also starts at A but *decreases* linearly with increasing B —its slope is $-\mu$. These two levels vary with B as shown in Fig. 12-3. We show also in the figure the energies E_{III} and E_{IV} . They have a different B -dependence. For small B , they depend quadratically on B , so they start out with horizontal slopes. Then they begin to curve, and for *large* B they approach straight lines with slopes $\pm\mu'$, which are nearly the same as the slopes of E_I and E_{II} .

The shift of the energy levels of an atom due to a magnetic field is called the *Zeeman effect*. We say that the curves in Fig. 12-3 show the *Zeeman splitting* of the ground state of hydrogen. When there is no magnetic field, we get just one spectral line from the hyperfine structure of hydrogen. The transitions between state $|IV\rangle$ and any one of the others occurs with the absorption or emission of a photon whose frequency 1420 megacycles is $1/h$ times the energy difference $4A$. When the atom is in a magnetic field B , however, there are many more lines. There can be transitions between any two of the four states. So if we have atoms in all four states, energy can be absorbed—or emitted—in any one of the six transitions shown by the vertical arrows in Fig. 12-4. Many of these transitions can be observed by the Rabi molecular beam technique we described in Volume II, Section 35-3 (see Appendix).

What makes the transitions go? The transitions will occur if you apply a small disturbing magnetic field that varies with time (in addition to the steady strong field B). It's just as we saw for a varying electric field on the ammonia molecule. Only here, it is the magnetic field which couples with the magnetic moments and does the trick. But the theory follows through in the same way that we worked it out for the ammonia. The theory is the simplest if you take a perturbing magnetic field that rotates in the xy -plane—although any horizontal oscillating field will do. When you put in this perturbing field as an additional term in the Hamiltonian, you get solutions in which the amplitudes vary with time—as we found for the ammonia molecule. So you can calculate easily and accurately the probability of a transition from one state to another. And you find that it all agrees with experiment.

12-5 The states in a magnetic field

We would like now to discuss the shapes of the curves in Fig. 12-3. In the first place, the energies for large fields are easy to understand, and rather interesting. For B large enough (namely for $\mu B/A \gg 1$) we can neglect the 1 in the formulas of (12.37). The four energies become

$$\begin{aligned} E_I &= A + \mu B, & E_{II} &= A - \mu B, \\ E_{III} &= -A + \mu' B, & E_{IV} &= -A - \mu' B. \end{aligned} \quad (12.38)$$

These are the equations of the four straight lines in Fig. 12-3. We can understand these energies physically in the following way. The nature of the stationary states in a *zero* field is determined completely by the interaction of the two magnetic moments. The mixtures of the base states $|+ -\rangle$ and $| - +\rangle$ in the stationary states $|III\rangle$ and $|IV\rangle$ are due to this interaction. In *large external* fields, however, the proton and electron will be influenced hardly at all by the field of the other; each will act as if it were alone in the external field. Then—as we have seen many times—the electron spin will be either parallel to or opposite to the external magnetic field.

Suppose the electron spin is “up”—that is, along the field; its energy will be $-\mu_e B$. The proton can still be either way. If the proton spin is also “up,” its energy is $-\mu_p B$. The sum of the two is $-(\mu_e + \mu_p)B = \mu B$. That is just what we find for E_I —which is fine, because we are describing the state $|+ +\rangle = |I\rangle$. There is still the small additional term A (now $\mu B \gg A$) which represents the interaction energy of the proton and electron when their spins are parallel. (We originally took A as positive because the theory we spoke of says it should be, and experimentally it is indeed so.) On the other hand, the proton can have its spin down. Then its energy in the external field goes to $-\mu_p B$, so it and the electron have the energy $-(\mu_e - \mu_p)B = \mu' B$. And the interaction energy becomes $-A$.

The sum is just the energy E_{III} in (12.38). So the state $|III\rangle$ must for large fields become the state $|+-\rangle$.

Suppose now the electron spin is “down.” Its energy in the external field is $\mu_e B$. If the proton is also “down,” the two together have the energy $(\mu_e + \mu_p)B = \mu B$, plus the interaction energy A —since their spins are parallel. That makes just the energy E_{II} in (12.38) and corresponds to the state $|--\rangle = |II\rangle$ —which is nice. Finally if the electron is “down” and the proton is “up,” we get the energy $(\mu_e - \mu_p)B - A$ (minus A for the interaction because the spins are opposite) which is just E_{IV} . And the state corresponds to $| - + \rangle$.

“But, wait a moment!”, you are probably saying, “The states $|III\rangle$ and $|IV\rangle$ are not the states $|+-\rangle$ and $| - + \rangle$; they are *mixtures* of the two.” Well, only slightly. They are indeed mixtures for $B = 0$, but we have not yet figured out what they are for large B . When we used the analogies of (12.33) in our formulas of Chapter 9 to get the energies of the stationary states, we could also have taken the amplitudes that go with them. They come from Eq. (9.23), which is

$$\frac{a_2}{a_3} = \frac{E - H_{22}}{H_{11}}.$$

The ratio a_2/a_3 is, of course, just C_2/C_3 . Plugging in the analogous quantities from (12.33), we get

$$\frac{C_2}{C_3} = \frac{E + A - (\mu_e - \mu_p)B}{2A}$$

or

$$\frac{C_2}{C_3} = \frac{E + A + \mu' B}{2A}, \quad (12.39)$$

where for E we are to use the appropriate energy—either E_{III} or E_{IV} . For instance, for state $|III\rangle$ we have

$$\left(\frac{C_2}{C_3}\right)_{III} \approx \frac{\mu' B}{A}. \quad (12.40)$$

So for large B the state $|III\rangle$ has $C_2 \gg C_3$; the state becomes almost completely the state $|2\rangle = |+-\rangle$. Similarly, if we put E_{IV} into (12.39) we get $(C_2/C_3)_{IV} \ll 1$; for high fields state $|IV\rangle$ becomes just the state $|3\rangle = |-+\rangle$. You see that the coefficients in the linear combinations of our base states which make up the stationary states depend on B . The state we call $|III\rangle$ is a 50–50 mixture of $|+-\rangle$ and $| - + \rangle$ at very low fields, but shifts completely over to $|+-\rangle$ at high fields. Similarly, the state $|IV\rangle$, which at low fields is also a 50–50 mixture (with opposite signs) of $|+-\rangle$ and $| - + \rangle$, goes over into the state $| - + \rangle$ when the spins are uncoupled by a strong external field.

We would also like to call your attention particularly to what happens at *very low* magnetic fields. There is one energy—at $-3A$ —which *does not change* when you turn on a small magnetic field. And there is another energy—at $+A$ —which splits into three different energy levels when you turn on a small magnetic field. For weak fields the energies vary with B as shown in Fig. 12-5. Suppose that we have somehow selected a bunch of hydrogen atoms which all have the energy $-3A$. If we put them through a Stern-Gerlach experiment—with fields that are not too strong—we would find that they just go straight through. (Since their energy doesn't depend on B , there is—according to the principle of virtual work—no force on them in a magnetic field gradient.) Suppose, on the other hand, we were to select a bunch of atoms with the energy $+A$, and put them through a Stern-Gerlach apparatus, say an S apparatus. (Again the fields in the apparatus should not be so great that they disrupt the insides of the atom, by which we mean a field small enough that the energies vary linearly with B .) We would find *three beams*. The states $|I\rangle$ and $|II\rangle$ get opposite forces—their energies vary linearly with B with the slopes $\pm\mu$ so the *forces* are like those on a dipole with $\mu_z = \mp\mu$; but the state $|III\rangle$ goes straight through. So we are right back in Chapter 5. A hydrogen atom with the energy $+A$ is a *spin-one particle*. This energy state is a “particle” for which $j = 1$, and it can be described—with respect to some set of

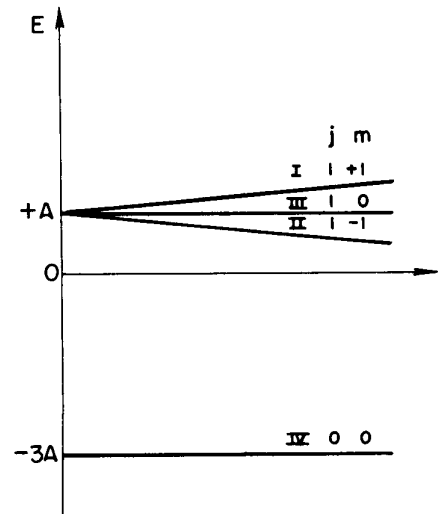


Fig. 12-5. The states of the hydrogen atom for small magnetic fields.

axes in space—in terms of the base states $|+S\rangle$, $|0S\rangle$, and $|-S\rangle$ we used in Chapter 5. On the other hand, when a hydrogen atom has the energy $-3A$, it is a spin-zero particle. (Remember, what we are saying is only strictly true for infinitesimal magnetic fields.) So we can group the states of hydrogen in zero magnetic field this way:

$$\left. \begin{aligned} |I\rangle &= |++\rangle \\ |III\rangle &= \frac{|+-\rangle + |-+\rangle}{\sqrt{2}} \\ |II\rangle &= |--\rangle \end{aligned} \right\} \text{spin } 1 \quad \left\{ \begin{aligned} &|+S\rangle \\ &|0S\rangle \\ &|-S\rangle \end{aligned} \right. \quad (12.41)$$

$$|IV\rangle = \frac{|+-\rangle - |-+\rangle}{\sqrt{2}} \quad \text{spin } 0. \quad (12.42)$$

We have said in Chapter 35 of Volume II (Appendix) that for any particle its component of angular momentum along any axis can have only certain values always \hbar apart. The z -component of angular momentum J_z can be $j\hbar$, $(j-1)\hbar$, $(j-2)\hbar$, \dots , $(-j)\hbar$, where j is the spin of the particle (which can be an integer or half-integer). Although we neglected to say so at the time, people usually write

$$J_z = m\hbar, \quad (12.43)$$

Table 12-3
Zero field states of the hydrogen atom

State $ j, m\rangle$	j	m	Our notation
$ 1, +1\rangle$	1	+1	$ I\rangle = +S\rangle$
$ 1, 0\rangle$	1	0	$ III\rangle = 0S\rangle$
$ 1, -1\rangle$	1	-1	$ II\rangle = -S\rangle$
$ 0, 0\rangle$	0	0	$ IV\rangle$

where m stands for one of the numbers $j, j-1, j-2, \dots, -j$. You will, therefore, see people in books label the four ground states of hydrogen by the so-called *quantum numbers* j and m [often called the “total angular momentum quantum number” (j), and “magnetic quantum number” (m)]. Then, instead of our state symbols $|I\rangle$, $|II\rangle$, and so on, they will write a state as $|j, m\rangle$. So they would write our little table of states for zero field in (12.41) and (12.42) as shown in Table 12-3. It’s not new physics, it’s all just a matter of notation.

12-6 The projection matrix for spin one †

We would like now to use our knowledge of the hydrogen atom to do something special. We discussed in Chapter 5 that a particle of *spin one* which was in one of the base states ($+$, 0 , or $-$) with respect to a Stern-Gerlach apparatus of a particular orientation—say an S apparatus—would have a certain amplitude to be in each of the three states with respect to a T apparatus with a different orientation in space. There are nine such amplitudes $\langle jT | iS \rangle$ which make up the projection matrix. In Section 5-7 we gave without proof the terms of this matrix for various orientations of T with respect to S . Now we will show you one way they can be derived.

In the hydrogen atom we have found a spin-one system which is made up of two spin one-half particles. We have already worked out in Chapter 6 how to transform the spin one-half amplitudes. We can use this information to calculate the transformation for spin one. This is the way it works: We have a system—a hydrogen atom with the energy $+A$ —which has spin one. Suppose we run it through a Stern-Gerlach filter S , so that we know it is in one of the base states with respect to S , say $|+S\rangle$. What is the amplitude that it will be in one of the base states, say $|+T\rangle$, with respect to the T apparatus? If we call the coordinate system of the S apparatus the x, y, z system, the $|+S\rangle$ state is what we have been calling the state $|++\rangle$. But suppose another guy took his z -axis along the axis of T . He will be referring his states to what we will call the x', y', z' frame. His “up” and “down” states for the electron and proton would be different from ours. His “plus-plus” state—which we can write $|++'\rangle$, referring to the “prime” frame—is the $|+T\rangle$ state of the spin-one particle. What we want is $\langle +T | +S \rangle$ which is just another way of writing the amplitude $\langle ++' | ++ \rangle$.

† Those who chose to jump over Chapter 6 should skip this section also.

We can find the amplitude $\langle +' + ' | + + \rangle$ in the following way. In *our* frame the *electron* in the $| + + \rangle$ state has its spin “up”. That means that it has some amplitude $\langle +' | + \rangle_e$ of being “up” in *his* frame, and some amplitude $\langle -' | + \rangle_e$ of being “down” in that frame. Similarly, the *proton* in the $| + + \rangle$ state has spin “up” in our frame and the amplitudes $\langle +' | + \rangle_p$ and $\langle -' | + \rangle_p$ of having spin “up” or spin “down” in the “prime” frame. Since we are talking about two distinct particles, the amplitude that *both* particles will be “up” *together* in *his* frame is the product of the two amplitudes,

$$\langle +' + ' | + + \rangle = \langle +' | + \rangle_e \langle +' | + \rangle_p. \quad (12.44)$$

We have put the subscripts *e* and *p* on the amplitudes $\langle +' | + \rangle$ to make it clear what we were doing. But they are both just the transformation amplitudes for a spin one-half particle, so they are really identical numbers. They are, in fact, just the amplitude we have called $\langle +T | +S \rangle$ in Chapter 6, and which we listed in the tables at the end of that chapter.

Now, however, we are about to get into trouble with notation. We have to be able to distinguish the amplitude $\langle +T | +S \rangle$ for a *spin one-half* particle from what we have *also* called $\langle +T | +S \rangle$ for a *spin-one* particle—yet they are completely different! We hope it won't be too confusing, but *for the moment* at least, we will have to use some different symbols for the spin one-half amplitudes. To help you keep things straight, we summarize the new notation in Table 12-4. We will continue to use the notation $| +S \rangle$, $| 0S \rangle$, and $| -S \rangle$ for the states of a spin-one particle.

With our new notation, Eq. (12.44) becomes simply

$$\langle +' + ' | + + \rangle = a^2,$$

and this is just the spin-one amplitude $\langle +T | +S \rangle$. Now, let's suppose, for instance, that the other guy's coordinate frame—that is, the *T*, or “primed,” apparatus—is just rotated with respect to *our* *z*-axis by the angle ϕ ; then from Table 6-2,

$$a = \langle +' | + \rangle = e^{i\phi/2}.$$

So from (12.44) we have that the spin-one amplitude is

$$\langle +T | +S \rangle = \langle +' + ' | + + \rangle = (e^{i\phi/2})^2 = e^{i\phi}. \quad (12.45)$$

You can see how it goes.

Now we will work through the general case for all the states. If the proton and electron are both “up” in *our* frame—the *S*-frame—the amplitudes that it will be in *any one of the four* possible states in the other guy's frame—the *T*-frame—are

$$\begin{aligned} \langle +' + ' | + + \rangle &= \langle +' | + \rangle_e \langle +' | + \rangle_p = a^2, \\ \langle +' -' | + + \rangle &= \langle +' | + \rangle_e \langle -' | + \rangle_p = ab, \\ \langle -' + ' | + + \rangle &= \langle -' | + \rangle_e \langle +' | + \rangle_p = ba, \\ \langle -' -' | + + \rangle &= \langle -' | + \rangle_e \langle -' | + \rangle_p = b^2. \end{aligned} \quad (12.46)$$

We can, then, write the state $| + + \rangle$ as the following linear combination:

$$| + + \rangle = a^2 | +' +' \rangle + ab \{ | +' -' \rangle + | -' +' \rangle \} + b^2 | -' -' \rangle. \quad (12.47)$$

Now we notice that $| +' +' \rangle$ is the state $| +T \rangle$, that $\{ | +' -' \rangle + | -' +' \rangle \}$ is just $\sqrt{2}$ times the state $| 0T \rangle$ —see (12.41)—and that $| -' -' \rangle = | -T \rangle$. In other words, Eq. (12.47) can be rewritten as

$$| +S \rangle = a^2 | +T \rangle + \sqrt{2} ab | 0T \rangle + b^2 | -T \rangle. \quad (12.48)$$

In a similar way you can easily show that

$$| -S \rangle = c^2 | +T \rangle + \sqrt{2} cd | 0T \rangle + d^2 | -T \rangle. \quad (12.49)$$

Table 12-4
Spin one-half amplitudes

This chapter	Chapter 6
$a = \langle +' + \rangle$	$\langle +T +S \rangle$
$b = \langle -' + \rangle$	$\langle -T +S \rangle$
$c = \langle +' - \rangle$	$\langle +T -S \rangle$
$d = \langle -' - \rangle$	$\langle -T -S \rangle$

For $|0 S\rangle$ it's a little more complicated, because

$$|0 S\rangle = \frac{1}{\sqrt{2}} \{|+-\rangle + |-+\rangle\}.$$

But we can express each of the states $|+-\rangle$ and $|-+\rangle$ in terms of the “prime” states and take the sum. That is,

$$|+-\rangle = ac|+'+\rangle + ad|+'-\rangle + bc|-' +'\rangle + bd|-' -'\rangle \quad (12.50)$$

and

$$|-+\rangle = ac|+'+\rangle + bc|+'-\rangle + ad|-' +'\rangle + bd|-' -'\rangle. \quad (12.51)$$

Taking $1/\sqrt{2}$ times the sum, we get

$$|0 S\rangle = \frac{2}{\sqrt{2}} ac|+'+\rangle + \frac{ad+bc}{\sqrt{2}} \{|+'-\rangle + |-' +'\rangle\} + \frac{2}{\sqrt{2}} bd|-' -'\rangle.$$

It follows that

$$|0 S\rangle = \sqrt{2} ac|+T\rangle + (ad+bc)|0 T\rangle + \sqrt{2} bd|-T\rangle. \quad (12.52)$$

We have now all of the amplitudes we wanted. The coefficients of Eqs. (12.48), (12.49), and (12.52) are the matrix elements $\langle jT|iS\rangle$. Let's pull them all together:

$$\langle jT|iS\rangle = \begin{matrix} iS \\ jT \downarrow \end{matrix} \begin{pmatrix} a^2 & \sqrt{2} ac & c^2 \\ \sqrt{2} ab & ad+bc & \sqrt{2} cd \\ b^2 & \sqrt{2} bd & d^2 \end{pmatrix} \quad (12.53)$$

We have expressed the spin-one transformation in terms of the spin one-half amplitudes a , b , c , and d .

For instance, if the T -frame is rotated with respect to S by the angle α about the y -axis—as in Fig. 5-6—the amplitudes in Table 12-4 are just the matrix elements of $R_y(\alpha)$ in Table 6-2.

$$\begin{aligned} a &= \cos \frac{\alpha}{2}, & b &= -\sin \frac{\alpha}{2}, \\ c &= \sin \frac{\alpha}{2}, & d &= \cos \frac{\alpha}{2}. \end{aligned} \quad (12.54)$$

Using these in (12.53), we get the formulas of (5.38), which we gave there without proof.

What ever happened to the state $|IV\rangle$?! Well, it is a spin-zero system, so it has only one state—it is the *same in all coordinate systems*. We can check that everything works out by taking the difference of Eq. (12.50) and (12.51); we get that

$$|+-\rangle - |-+\rangle = (ad-bc)\{|+'-\rangle - |-' +'\rangle\}.$$

But $(ad-bc)$ is the determinant of the spin one-half matrix, and so is equal to 1. We get that

$$|IV'\rangle = |IV\rangle$$

for any relative orientation of the two coordinate frames.

Propagation in a Crystal Lattice

13-1 States for an electron in a one-dimensional lattice

You would, at first sight, think that a low-energy electron would have great difficulty passing through a solid crystal. The atoms are packed together with their centers only a few angstroms apart, and the effective diameter of the atom for electron scattering is roughly an angstrom or so. That is, the atoms are large, relative to their spacing, so that you would expect the mean free path between collisions to be of the order of a few angstroms—which is practically nothing. You would expect the electron to bump into one atom or another almost immediately. Nevertheless, it is a ubiquitous phenomenon of nature that if the lattice is perfect, the electrons are able to travel through the crystal smoothly and easily—almost as if they were in a vacuum. This strange fact is what lets metals conduct electricity so easily; it has also permitted the development of many practical devices. It is, for instance, what makes it possible for a transistor to imitate the radio tube. In a radio tube electrons move freely through a vacuum, while in the transistor they move freely through a crystal lattice. The machinery behind the behavior of a transistor will be described in this chapter; the next one will describe the application of these principles in various practical devices.

The conduction of electrons in a crystal is one example of a very common phenomenon. Not only can electrons travel through crystals, but other “things” like atomic excitations can also travel in a similar manner. So the phenomenon which we want to discuss appears in many ways in the study of the physics of the solid state.

You will remember that we have discussed many examples of two-state systems. Let’s now think of an electron which can be in either one of two positions, in each of which it is in the same kind of environment. Let’s also suppose that there is a certain amplitude to go from one position to the other, and, of course, the same amplitude to go back, just as we have discussed for the hydrogen molecular ion in Section 10-1. The laws of quantum mechanics then give the following results. There are two possible states of definite energy for the electron. Each state can be described by the amplitude for the electron to be in each of the two basic positions. In either of the definite-energy states, the magnitudes of these two amplitudes are constant in time, and the phases vary in time with the same frequency. On the other hand, if we start the electron in one position, it will later have moved to the other, and still later will swing back again to the first position. The amplitude is analogous to the motions of two coupled pendulums.

Now consider a perfect crystal lattice in which we imagine that an electron can be situated in a kind of “pit” at one particular atom and with some particular energy. Suppose also that the electron has some amplitude to move into a different pit at one of the nearby atoms. It is something like the two-state system—but with an additional complication. When the electron arrives at the neighboring atom, it can afterward move on to still another position as well as return to its starting point. Now we have a situation analogous not to *two* coupled pendulums, but to an *infinite number* of pendulums all coupled together. It is something like what you see in one of those machines—made with a long row of bars mounted on a torsion wire—that is used in first-year physics to demonstrate wave propagation.

If you have a harmonic oscillator which is coupled to another harmonic oscillator, and that one to another, and so on . . . , and if you start an irregularity in one place, the irregularity will propagate as a wave along the line. The same situation exists if you place an electron at one atom of a long chain of atoms.

13-1 States for an electron in a one-dimensional lattice

13-2 States of definite energy

13-3 Time-dependent states

13-4 An electron in a three-dimensional lattice

13-5 Other states in a lattice

13-6 Scattering by imperfections in the lattice

13-7 Trapping by a lattice imperfection

13-8 Scattering amplitudes and bound states

Usually, the simplest way of analyzing the mechanical problem is not to think in terms of what happens if a pulse is started at a definite place, but rather in terms of steady-wave solutions. There exist certain patterns of displacements which propagate through the crystal as a wave of a single, fixed frequency. Now the same thing happens with the electron—and for the same reason, because it's described in quantum mechanics by similar equations.

You must appreciate one thing, however; the amplitude for the electron to be at a place is an *amplitude*, not a probability. If the electron were simply leaking from one place to another, like water going through a hole, the behavior would be completely different. For example, if we had two tanks of water connected by a tube to permit some leakage from one to the other, then the levels would approach each other exponentially. But for the electron, what happens is amplitude leakage and not just a plain probability leakage. And it's a characteristic of the imaginary term—the i in the differential equations of quantum mechanics—which changes the exponential solution to an oscillatory solution. What happens then is quite different from the leakage between interconnected tanks.

We want now to analyze quantitatively the quantum mechanical situation. Imagine a one-dimensional system made of a long line of atoms as shown in Fig. 13-1(a). (A crystal is, of course, three-dimensional but the physics is very much the same; once you understand the one-dimensional case you will be able to understand what happens in three dimensions.) Next, we want to see what happens if we put a single electron on this line of atoms. Of course, in a real crystal there are already millions of electrons. But most of them (nearly all for an insulating crystal) take up positions in some pattern of motion each around its own atom—and everything is quite stationary. However, we now want to think about what happens if we put an *extra* electron in. We will not consider what the other ones are doing because we suppose that to change their motion involves a lot of excitation energy. We are going to add an electron as if to produce one slightly bound negative ion. In watching what the *one* extra electron does we are making an approximation which disregards the mechanics of the inside workings of the atoms.

Of course the electron could then move to another atom, transferring the negative ion to another place. We will suppose that just as in the case of an electron jumping between two protons, the electron can jump from one atom to the neighbor on either side with a certain amplitude.

Now how do we describe such a system? What will be reasonable base states? If you remember what we did when we had only two possible positions, you can guess how it will go. Suppose that in our line of atoms the spacings are all equal; and that we number the atoms in sequence, as shown in Fig. 13-1(a). One of the base states is that the electron is at atom number 6, another base state is that the electron is at atom number 7, or at atom number 8, and so on. We can describe the n th base state by saying that the electron is at atom number n . Let's say that this is the base state $|n\rangle$. Figure 13-1 shows what we mean by the three base states

$$|n-1\rangle, |n\rangle, \text{ and } |n+1\rangle.$$

Using these base states, any state $|\phi\rangle$ of our one-dimensional crystal can be described by giving all the amplitudes $\langle n|\phi\rangle$ that the state $|\phi\rangle$ is in one of the base states—which means the amplitude that it is located at one particular atom. Then we can write the state $|\phi\rangle$ as a superposition of the base states

$$|\phi\rangle = \sum_n |n\rangle \langle n|\phi\rangle. \quad (13.1)$$

Next, we are going to suppose that when the electron is at one atom, there is a certain amplitude that it will leak to the atom on either side. And we'll take the simplest case for which it can only leak to the nearest neighbors—to get to the next-nearest neighbor, it has to go in two steps. We'll take that the amplitudes for the electron jump from one atom to the next is iA/\hbar (per unit time).

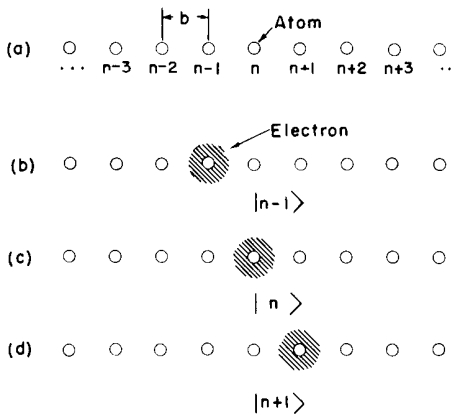


Fig. 13-1. The base states of an electron in a one-dimensional crystal.

For the moment we would like to write the amplitude $\langle n | \phi \rangle$ to be on the n th atom as C_n . Then Eq. (13.1) will be written

$$|\phi\rangle = \sum_n |n\rangle C_n. \quad (13.2)$$

If we knew each of the amplitudes C_n at a given moment, we could take their absolute squares and get the probability that you would find the electron if you looked at atom n at that time.

What will the situation be at some later time? By analogy with the two-state systems we have studied, we would propose that the Hamiltonian equations for this system should be made up of equations like this:

$$i\hbar \frac{dC_n(t)}{dt} = E_0 C_n(t) - A C_{n+1}(t) - A C_{n-1}(t). \quad (13.3)$$

The first coefficient on the right, E_0 , is, physically, the energy the electron would have if it couldn't leak away from one of the atoms. (It doesn't matter what we call E_0 ; as we have seen many times, it represents really nothing but our choice of the zero of energy.) The next term represents the amplitude per unit time that the electron is leaking into the n th pit from the $(n + 1)$ st pit; and the last term is the amplitude for leakage from the $(n - 1)$ st pit. As usual, we'll assume that A is a constant (independent of t).

For a full description of the behavior of any state $|\phi\rangle$, we would have one equation like (13.3) for every one of the amplitudes C_n . Since we want to consider a crystal with a very large number of atoms, we'll assume that there are an indefinitely large number of states—that the atoms go on forever in both directions. (To do the finite case, we will have to pay special attention to what happens at the ends.) If the number N of our base states is indefinitely large, then also our full Hamiltonian equations are infinite in number! We'll write down just a sample:

$$\begin{array}{ccc} \vdots & & \vdots \\ i\hbar \frac{dC_{n-1}}{dt} & = & E_0 C_{n-1} - A C_{n-2} - A C_n, \\ i\hbar \frac{dC_n}{dt} & = & E_0 C_n - A C_{n-1} - A C_{n+1}, \\ i\hbar \frac{dC_{n+1}}{dt} & = & E_0 C_{n+1} - A C_n - A C_{n+2}, \\ \vdots & & \vdots \end{array} \quad (13.4)$$

13-2 States of definite energy

We could study many things about an electron in a lattice, but first let's try to find the states of definite energy. As we have seen in earlier chapters this means that we have to find a situation in which the amplitudes all change at the same frequency if they change with time at all. We look for solutions of the form

$$C_n = a_n e^{-iE t/\hbar}. \quad (13.5)$$

The complex number a_n tell us about the non-time-varying part of the amplitude to find the electron at the n th atom. If we put this trial solution into the equations of (13.4) to test them out, we get the result

$$E a_n = E_0 a_n - A a_{n+1} - A a_{n-1}. \quad (13.6)$$

We have an infinite number of such equations for the infinite number of unknowns a_n —which is rather petrifying.

All we have to do is take the determinant . . . but wait! Determinants are fine when there are 2, 3, or 4 equations. But if there are a large number—or an infinite number—of equations, the determinants are not very convenient. We'd better just try to solve the equations directly. First, let's label the atoms by their

positions; we'll say that the atom n is at x_n and the atom $(n + 1)$ is at x_{n+1} . If the atomic spacing is b —as in Fig. 13-1—we will have that $x_{n+1} = x_n + b$. By choosing our origin at atom zero, we can even have it that $x_n = nb$. We can rewrite Eq. (13.5) as

$$C_n = a(x_n)e^{-iEt/\hbar}, \quad (13.7)$$

and Eq. (13.6) would become

$$Ea(x_n) = E_0a(x_{n+1}) - Aa(x_{n+1}) - Aa(x_{n-1}). \quad (13.8)$$

Or, using the fact that $x_{n+1} = x_n + b$, we could also write

$$Ea(x_n) = E_0a(x_n) - Aa(x_n + b) - Aa(x_n - b). \quad (13.9)$$

This equation is somewhat similar to a differential equation. It tells us that a quantity, $a(x)$, at one point, (x_n) , is related to the same physical quantity at some neighboring points, $(x_n \pm b)$. (A differential equation relates the value of a function at a point to the values at infinitesimally nearby points.) Perhaps the methods we usually use for solving differential equations will also work here, let's try.

Linear differential equations with constant coefficients can always be solved in terms of exponential functions. We can try the same thing here; let's take as a trial solution

$$a(x_n) = e^{ikx_n}. \quad (13.10)$$

Then Eq. (13.9) becomes

$$Ee^{ikx_n} = E_0e^{ikx_n} - Ae^{ik(x_n+b)} - Ae^{ik(x_n-b)}. \quad (13.11)$$

We can now divide out the common factor e^{ikx_n} ; we get

$$E = E_0 - Ae^{ikb} - Ae^{-ikb}. \quad (13.12)$$

The last two terms are just equal to $(2A \cos kb)$, so

$$E = E_0 - 2A \cos kb. \quad (13.13)$$

We have found that for *any* choice at all for the constant k there is a solution whose energy is given by this equation. There are various possible energies depending on k , and each k corresponds to a different solution. There are an infinite number of solutions—which is not surprising, since we started out with an infinite number of base states.

Let's see what these solutions mean. For each k , the a 's are given by Eq. (13.10). The amplitudes C_n are then given by

$$C_n = e^{ikx_n}e^{-(i/\hbar)Et}, \quad (13.14)$$

where you should remember that the energy E also depends on k as given in Eq. (13.13). The *space dependence* of the amplitudes is e^{ikx_n} . The amplitudes oscillate as we go along from one atom to the next.

We mean that, in space, the amplitude goes as a *complex* oscillation—the *magnitude* is the same at every atom, but the phase at a given time advances by the amount (ikb) from one atom to the next. We can visualize what is going on by plotting a vertical line to show just the real part at each atom as we have done in Fig. 13-2. The envelope of these vertical lines (as shown by the broken-line curve)

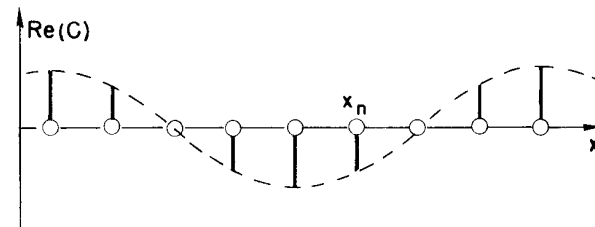


Fig. 13-2. Variation of the real part of C_n with x_n .

is, of course, a cosine curve. The imaginary part of C_n is also an oscillating function, but is shifted 90° in phase so that the absolute square (which is the sum of the squares of the real and imaginary parts) is the same for all the C 's.

Thus if we pick a k , we get a stationary state of a particular energy E . And for any such state, the electron is equally likely to be found at every atom—there is no preference for one atom or the other. Only the phase is different for different atoms. Also, as time goes on the phases vary. From Eq. (13.14) the real and imaginary parts propagate along the crystal as waves—namely as the real or imaginary parts of

$$e^{i[kx_n - (E/\hbar)t]}. \quad (13.15)$$

The wave can travel toward positive or negative x depending on the sign we have picked for k .

Notice that we have been assuming that the number k that we put in our trial solution, Eq. (13.10), was a real number. We can see now why that must be so if we have an infinite line of atoms. Suppose that k were an imaginary number, say ik' . Then the amplitudes a_n would go as $e^{k'x_n}$, which means that the amplitude would get larger and larger as we go toward large x 's—or toward large negative x 's if k' is a negative number. This kind of solution would be O.K. if we were dealing with line of atoms that ended, but cannot be a physical solution for an infinite chain of atoms. It would give infinite amplitudes—and, therefore, infinite probabilities—which can't represent a real situation. Later on we will see an example in which an imaginary k does make sense.

The relation between the energy E and the wave number k as given in Eq. (13.13) is plotted in Fig. 13-3. As you can see from the figure, the energy can go from $(E_0 - 2A)$ at $k = 0$ to $(E_0 + 2A)$ at $k = \pm\pi/b$. The graph is plotted for positive A ; if A were negative, the curve would simply be inverted, but the range would be the same. The significant result is that any energy is possible within a certain range or "band" of energies, but no others. According to our assumptions, if an electron in a crystal is in a stationary state, it can have no energy other than values in this band.

According to Eq. (13.10), the smallest k 's correspond to low-energy states— $E \approx (E_0 - 2A)$. As k increases in magnitude (toward either positive or negative values) the energy at first increases, but then reaches a maximum at $k = \pm\pi/b$, as shown in Fig. 13-3. For k 's larger than π/b , the energy would start to decrease again. But we do not really need to consider such values of k , because they do not give new states—they just repeat states we already have for smaller k . We can see that in the following way. Consider the lowest energy state for which $k = 0$. The coefficient $a(x_n)$ is the same for all x_n . Now we would get the same energy for $k = 2\pi/b$. But then, using Eq. (13.10), we have that

$$a(x_n) = e^{i(2\pi/b)x_n}.$$

However, taking x_0 to be at the origin, we can set $x_n = nb$; then $a(x_n)$ becomes

$$a(x_n) = e^{i2\pi n} = 1.$$

The state described by these $a(x_n)$ is physically the same state we got for $k = 0$. It does not represent a different solution.

As another example, suppose that k were $\pi/4b$. The real part of $a(x_n)$ would vary as shown by curve 1 in Fig. 13-4. If k were seven times larger ($k = 7\pi/4$), the real part of $a(x_n)$ would vary as shown by curve 2 in the figure. (The complete

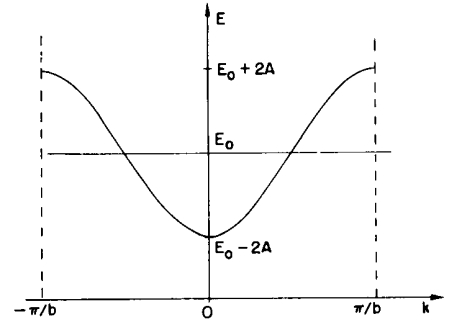


Fig. 13-3. The energy of the stationary states as a function of the parameter k .

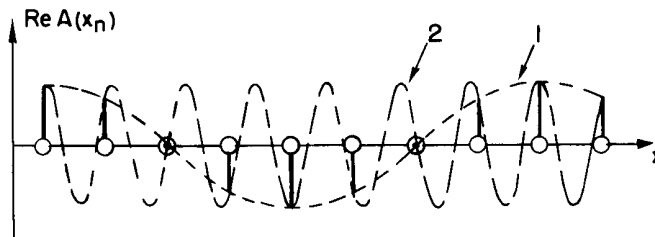


Fig. 13-4. Two values of k which represent the same physical situation; curve 1 is for $k = \pi/4$, curve 2 is for $k = 7\pi/4$.

cosine curves don't mean anything, of course; all that matters is their values *at the points* x_n . The curves are just to help you see how things are going.) You see that both values of k give the same amplitudes at all of the x_n 's.

The upshot is that we have all the possible solutions of our problem if we take only k 's in a certain limited range. We'll pick the range between $-\pi/b$ and $+\pi/b$ —the one shown in Fig. 13-3. In this range, the energy of the stationary states increases uniformly with an increase in the magnitude of k .

One side remark about something you can play with. Suppose that the electron cannot only jump to the nearest neighbor with amplitude iA/\hbar , but also has the possibility to jump in one direct leap to the *next nearest* neighbor with some other amplitude iB/\hbar . You will find that the solution can again be written in the form $a_n = e^{ikx_n}$ —this type of solution is universal. You will also find that the stationary states with wave number k have an energy equal to $(E_0 - 2A \cos kb - 2B \cos 2kb)$. This shows that the shape of the curve of E against k is not universal, but depends upon the particular assumptions of the problem. It is not always a cosine wave—it's not even necessarily symmetrical about some horizontal line. It is true, however, that the curve always repeats itself outside of the interval from $-\pi/b$ to π/b , so you never need to worry about other values of k .

Let's look a little more closely at what happens for small k —that is, when the variations of the amplitudes from one x_n to the next are quite slow. Suppose we choose our zero of energy by defining $E_0 = 2A$; then the minimum of the curve in Fig. 13-3 is at the zero of energy. For small enough k , we can write that

$$\cos kb \approx 1 - k^2 b^2 / 2,$$

and the energy of Eq. (13.13) becomes

$$E = Ak^2 b^2. \quad (13.16)$$

We have that the energy of the state is proportional to the square of the wave number which describes the spatial variations of the amplitudes C_n .

13-3 Time-dependent states

In this section we would like to discuss the behavior of states in the one-dimensional lattice in more detail. If the amplitude for an electron to be at x_n is C_n , the probability of finding it there is $|C_n|^2$. For the *stationary* states described by Eq. (13.12), this probability is the same for all x_n and does not change with time. How can we represent a situation which we would describe roughly by saying an electron of a certain energy is localized in a certain region—so that it is more likely to be found at one place than at some other place? We can do that by making a superposition of several solutions like Eq. (13.12) with slightly different values of k —and, therefore, slightly different energies. Then at $t = 0$, at least, the amplitude C_n will vary with position because of the interference between the various terms, just as one gets beats when there is a mixture of waves of different wavelengths (as we discussed in Chapter 48, Vol. I). So we can make up a “wave packet” with a predominant wave number k_0 , but with various other wave numbers near k_0 .†

In our superposition of stationary states, the amplitudes with different k 's will represent states of slightly different energies, and, therefore, of slightly different frequencies; the interference pattern of the total C_n will, therefore, also vary with time—there will be a pattern of “beats.” As we have seen in Chapter 48 of Volume I, the peaks of the beats [the place where $|C(x_n)|^2$ is large] will move along in x as time goes on; they move with the speed we have called the “group velocity.” We found that this group velocity was related to the variation of k with frequency by

$$v_{\text{group}} = \frac{d\omega}{dk}; \quad (13.17)$$

† Provided we do not try to make the packet too narrow.

the same derivation would apply equally well here. An electron state which is a “clump”—namely one for which the C_n vary in space like the wave packet of Fig. 13-5—will move along our one-dimensional “crystal” with the speed v equal to $d\omega/dk$, where $\omega = E/\hbar$. Using (13.16) for E , we get that

$$v = \frac{2Ab^2}{\hbar} k. \quad (13.18)$$

In other words, the electrons move along with a speed proportional to the typical k . Equation (13.16) then says that the energy of such an electron is proportional to the square of its velocity—it acts like a classical particle. So long as we look at things on a scale gross enough that we don't see the fine structure, our quantum mechanical picture begins to give results like classical physics. In fact, if we solve Eq. (13.18) for k and substitute into (13.16), we can write

$$E = \frac{1}{2}m_{\text{eff}} v^2, \quad (13.19)$$

where m_{eff} is a constant. The extra “energy of motion” of the electron in a packet depends on the velocity just as for a classical particle. The constant m_{eff} —called the “effective mass”—is given by

$$m_{\text{eff}} = \frac{\hbar^2}{2Ab^2}. \quad (13.20)$$

Also notice that we can write

$$m_{\text{eff}} v = \hbar k. \quad (13.21)$$

If we choose to call $m_{\text{eff}} v$ the “momentum,” it is related to the wave number k in the way we have described earlier for a free particle.

Don't forget that m_{eff} has nothing to do with the real mass of an electron. It may be quite different—although in real crystals it often happens to turn out to be the same general order of magnitude, about 2 to 20 times the free-space mass of the electron.

We have now explained a remarkable mystery—how an electron in a crystal (like an extra electron put into germanium) can ride right through the crystal and flow perfectly freely even though it has to hit all the atoms. It does so by having its amplitudes going pip-pip-pip from one atom to the next, working its way through the crystal. That is how a solid can conduct electricity.

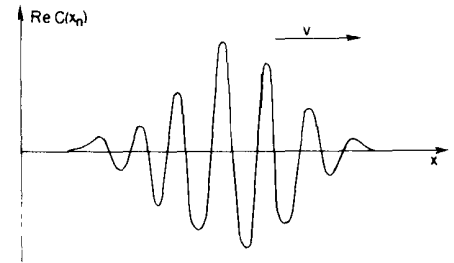


Fig. 13-5. The real part of $C(x_n)$ as a function of x for a superposition of several states of similar energy. (The spacing b is very small on the scale of x shown.)

13-4 An electron in a three-dimensional lattice

Let's look for a moment at how we could apply the same ideas to see what happens to an electron in three dimensions. The results turn out to be very similar. Suppose we have a rectangular lattice of atoms with lattice spacings of a , b , c in the three directions (If you want a cubic lattice, take the three spacings all equal.) Also suppose that the amplitude to leap in the x -direction to a neighbor is (tA_x/\hbar) , to leap in the y -direction is (tA_y/\hbar) , and to leap in the z -direction is (tA_z/\hbar) . Now how should we describe the base states? As in the one-dimensional case, one base state is that the electron is at the atom whose locations are x , y , z , where (x, y, z) is one of the lattice points. Choosing our origin at one atom, these points are all at

$$x = n_x a, \quad y = n_y b, \quad \text{and} \quad z = n_z c,$$

where n_x , n_y , n_z are any three integers. Instead of using subscripts to indicate such points, we will now just use x , y , and z , understanding that they take on only their values at the lattice points. Thus the base state is represented by the symbol |electron at x, y, z ⟩, and the amplitude for an electron in some state $|\psi\rangle$ to be in this base state is $C(x, y, z) = \langle \text{electron at } x, y, z | \psi \rangle$.

As before, the amplitudes $C(x, y, z)$ may vary with time. With our assumptions, the Hamiltonian equations should be like this:

$$i\hbar \frac{dC(x, y, z)}{dt} = E_0 C(x, y, z) - A_x C(x + a, y, z) - A_x C(x - a, y, z) \\ - A_y C(x, y + b, z) - A_y C(x, y - b, z) \\ - A_z C(x, y, z + c) - A_z C(x, y, z - c). \quad (13.22)$$

It looks rather long, but you can see where each term comes from.

Again we can try to find a stationary state in which all the C 's vary with time in the same way. Again the solution is an exponential:

$$C(x, y, z) = e^{-iEt/\hbar} e^{i(k_x x + k_y y + k_z z)}. \quad (13.23)$$

If you substitute this into (13.22) you see that it works, provided that the energy E is related to k_x , k_y , and k_z in the following way:

$$E = E_0 - 2A_x \cos k_x a - 2A_y \cos k_y b - 2A_z \cos k_z c. \quad (13.24)$$

The energy now depends on the *three* wave numbers k_x , k_y , k_z , which, incidentally, are the components of a three-dimensional vector \mathbf{k} . In fact, we can write Eq. (13.23) in vector notation as

$$C(x, y, z) = e^{-iEt/\hbar} e^{-i\mathbf{k} \cdot \mathbf{r}} \quad (13.25)$$

The amplitude varies as a complex *plane wave* in three dimensions, moving in the direction of \mathbf{k} , and with the wave number $k = (k_x^2 + k_y^2 + k_z^2)^{1/2}$.

The energy associated with these stationary states depends on the three components of \mathbf{k} in the complicated way given in Eq. (13.24). The nature, of the variation of E with \mathbf{k} depends on relative signs and magnitudes of A_x , A_y , and A_z . If these three numbers are all positive, and if we are interested in small values of k , the dependence is relatively simple.

Expanding the cosines as we did before to get Eq. (13.16), we can now get that

$$E = E_{\min} + A_x a^2 k_x^2 + A_y b k_y^2 + A_z c k_z^2. \quad (13.26)$$

For a simple cubic lattice with lattice spacing a we expect that A_x and A_y and A_z would be equal—say all are just A —and we would have just

$$E = E_{\min} + A a^2 (k_x^2 + k_y^2 + k_z^2),$$

or

$$E = E_{\min} + A a^2 k^2. \quad (13.27)$$

This is just like Eq. (13.16). Following the arguments used there, we would conclude that an electron packet in *three* dimensions (made up by superposing many states with nearly equal energies) also moves like a classical particle with some effective mass.

In a crystal with a lower symmetry than cubic (or even in a cubic crystal in which the state of the electron at each atom is not symmetrical) the three coefficients A_x , A_y , and A_z are different. Then the “effective mass” of an electron localized in a small region *depends on its direction of motion*. It could, for instance, have a different inertia for motion in the x -direction than for motion in the y -direction. (The details of such a situation are sometimes described in terms of an “effective mass tensor.”)

13-5 Other states in a lattice

According to Eq. (13.24) the electron states we have been talking about can have energies only in a certain “band” of energies which covers the energy range from the minimum energy

$$E_0 - 2(A_x + A_y + A_z)$$

to the maximum energy

$$E_0 + 2(A_x + A_y + A_z).$$

Other energies are possible, but they belong to a different class of electron states. For the states we have described, we imagined base states in which an electron is placed on an atom of the crystal in some particular state, say the lowest energy state.

If you have an atom in empty space, and add an electron to make an ion, the ion can be formed in many ways. The electron can go on in such a way as to make the state of lowest energy, or it can go on to make one or another of many possible “excited states” of the ion each with a definite energy above the lowest energy. The same thing can happen in a crystal. Let’s suppose that the energy E_0 we picked above corresponds to base states which are ions of the lowest possible energy. We could also imagine a new set of base states in which the electron sits near the n th atom in a different way—in one of the excited states of the ion—so that the energy E_0 is now quite a bit higher. As before there is some amplitude A (different from before) that the electron will jump from its excited state at one atom to the same excited state at a neighboring atom. The whole analysis goes as before, we find a band of possible energies centered at a higher energy. There can, in general, be many such bands each corresponding to a different level of excitation.

There are also other possibilities. There may be some amplitude that the electron jumps from an excited condition at one atom to an unexcited condition at the next atom. (This is called an interaction between bands.) The mathematical theory gets more and more complicated as you take into account more and more bands and add more and more coefficients for leakage between the possible states. No new ideas are involved, however; the equations are set up much as we have done in our simple example.

We should remark also that there is not much more to be said about the various coefficients, such as the amplitude A , which appear in the theory. Generally they are very hard to calculate, so in practical cases very little is known theoretically about these parameters and for any particular real situation we can only take values determined experimentally.

There are other situations where the physics and mathematics are almost exactly like what we have found for an electron moving in a crystal, but in which the “object” that moves is quite different. For instance, suppose that our original crystal—or rather linear lattice—was a line of neutral atoms, each with a loosely bound outer electron. Then imagine that we were to remove one electron. Which atom has lost its electron? Let C_n now represent the amplitude that the electron is *missing* from the atom at x_n . There will, in general, be some amplitude iA/\hbar that the electron at a neighboring atom—say the $(n - 1)$ st atom—will jump to the n th leaving the $(n - 1)$ st atom without its electron. This is the same as saying that there is an amplitude A for the “missing electron” to jump from the n th atom to the $(n - 1)$ st atom. You can see that the equations will be exactly the same—of course, the value of A need not be the same as we had before. Again we will get the same formulas for the energy levels, for the “waves” of probability which move through the crystal with the group velocity of Eq. (13.18), for the effective mass, and so on. Only now the waves describe the behavior of the *missing electron*—or “hole” as it is called. So a “hole” acts just like a particle with a certain mass m_{eff} . You can see that this particle will appear to have a positive charge. We’ll have some more to say about such holes in the next chapter.

As another example, we can think of a line of identical *neutral* atoms one of which has been put into an excited state—that is, with more than its normal ground state energy. Let C_n be the amplitude that the n th atom has the excitation. It can interact with a neighboring atom by handing over to it the extra energy and returning to the ground state. Call the amplitude for this process iA/\hbar . You can see that it’s the same mathematics all over again. Now the object which moves is called an *exciton*. It behaves like a neutral “particle” moving through the crystal, carrying the excitation energy. Such motion may be involved in certain biological

processes such as vision, or photosynthesis. It has been guessed that the absorption of light in the retina produces an “exciton” which moves through some periodic structure (such as the layers in the rods we described in Chapter 36, Vol. 1; see Fig. 36-5) to be accumulated at some special station where the energy is used to induce a chemical reaction.

13-6 Scattering from imperfections in the lattice

We want now to consider the case of a single electron in a crystal which is not perfect. Our earlier analysis says that perfect crystals have perfect conductivity—that electrons can go slipping through the crystal, as in a vacuum, without friction. One of the most important things that can stop an electron from going on forever is an imperfection or irregularity in the crystal. As an example, suppose that somewhere in the crystal there is a missing atom; or suppose that someone put one wrong atom at one of the atomic sites so that things there are different than at the other atomic sites. Say the energy, E_0 or the amplitude A could be different. How would we describe what happens then?

To be specific, we will return to the one-dimensional case and we will assume that atom number “zero” is an “impurity” atom and has a different value of E_0 than any of the other atoms. Let’s call this energy ($E_0 + F$). What happens? When an electron arrives at atom “zero” there is some probability that the electron is scattered backwards. If a wave packet is moving along and it reaches a place where things are a little bit different, some of it will continue onward and some of it will bounce back. It’s quite difficult to analyze such a situation using a wave packet, because everything varies in time. It is much easier to work with steady-state solutions. So we will work with stationary states, which we will find can be made up of continuous waves which have transmitted and reflected parts. In three dimensions we would call the reflected part the scattered wave, since it would spread out in various directions.

We start out with a set of equations which are just like the ones in Eq. (13.6) except that the equation for $n = 0$ is different from all the rest. The five equations for $n = -2, -1, 0, +1$, and $+2$ look like this:

$$\begin{aligned}
 & \vdots & & \vdots \\
 Ea_{-2} &= E_0 a_{-2} - Aa_{-1} - Aa_{-3}, \\
 Ea_{-1} &= E_0 a_{-1} - Aa_0 - Aa_{-2}, \\
 Ea_0 &= (E_0 + F)a_0 - Aa_1 - Aa_{-1}, \\
 Ea_1 &= E_0 a_1 - Aa_2 - Aa_0, \\
 Ea_2 &= E_0 a_2 - Aa_3 - Aa_1, \\
 & \vdots & & \vdots
 \end{aligned} \tag{13.28}$$

There are, of course, all the other equations for $|n|$ is greater than 2. They will look just like Eq. (13.16).

For the general case, we really ought to use a different A for the amplitude that the electron jumps to or from atom “zero,” but the main features of what goes on will come out of a simplified example in which all the A ’s are equal.

Equation (13.10) would still work as a solution for all of the equations except the one for atom “zero”—it isn’t right for that one equation. We need a different solution which we can cook up in the following way. Equation (13.10) represents a wave going in the positive x -direction. A wave going in the negative x -direction would have been an equally good solution. It would be written

$$a(x_n) = e^{-ikx_n}.$$

The most general solution we could have taken for Eq. (13.6) would be a com-

bination of a forward and a backward wave, namely

$$a_n = \alpha e^{ikx_n} + \beta e^{-ikx_n}. \quad (13.29)$$

This solution represents a complex wave of amplitude α moving in the $+x$ -direction and a wave of amplitude β moving in the $-x$ -direction.

Now take a look at the set of equations for our new problem—the ones in (13.28) together with those for all the other atoms. The equations involving a_n 's with $n \leq 1$ are all satisfied by Eq. (13.29), with the condition that k is related to E and the lattice spacing b by

$$E = E_0 - 2A \cos kb. \quad (13.30)$$

The physical meaning is an “incident” wave of amplitude α approaching atom “zero” (the “scatterer”) from the left, and a “scattered” or “reflected” wave of amplitude β going back toward the left. We do not lose any generality if we set the amplitude α of the incident wave equal to 1. Then the amplitude β is, in general, a complex number.

We can say all the same things about the solutions of a_n for $n \geq 1$. The coefficients could be different, so we would have for them

$$a_n = \gamma e^{ikx_n} + \delta e^{-ikx_n}, \quad \text{for } n \geq 1. \quad (13.31)$$

Here, γ is the amplitude of a wave going to the right and δ a wave coming from the right. We want to consider the *physical* situation in which a wave is originally started only from the left, and there is only a “transmitted” wave that comes out beyond the scatterer—or impurity atom. We will try for a solution in which $\delta = 0$. We can, certainly, satisfy all of the equations for the a_n except for the middle three in Eq. (13.28) by the following trial solutions.

$$\begin{aligned} a_n \text{ (for } n < 0) &= e^{ikx_n} + \beta e^{-ikx_n}, \\ a_n \text{ (for } n > 0) &= \gamma e^{ikx_n}. \end{aligned} \quad (13.32)$$

The situation we are talking about is illustrated in Fig. 13-6.

By using the formulas in Eq. (13.32) for a_{-1} and a_{+1} , the three middle equations of Eq. (13.28) will allow us to solve for a_0 and also for the two coefficients β and γ . So we have found a complete solution. Setting $x_n = nb$, we have to solve the three equations

$$\begin{aligned} (E - E_0)\{e^{ik(-b)} + \beta e^{-ik(-b)}\} &= -A\{a_0 + e^{ik(-2b)} + \beta e^{-ik(-2b)}\}, \\ (E - E_0 - F)a_0 &= -A\{\gamma e^{ikb} + e^{ik(-b)} + \beta e^{-ik(-b)}\}, \\ (E - E_0)\gamma e^{ikb} &= -A\{\gamma e^{ik(2b)} + a_0\}. \end{aligned} \quad (13.33)$$

Remember that E is given in terms of k by Eq. (13.30). If you substitute this value for E into the equations, and remember that $\cos x = \frac{1}{2}(e^{ix} + e^{-ix})$, you get from the first equation that

$$a_0 = 1 + \beta; \quad (13.34)$$

and from the third equation that

$$a_0 = \gamma. \quad (13.35)$$

These are consistent only if

$$\gamma = 1 + \beta \quad (13.36)$$

This equation says that the transmitted wave (γ) is just the original incident wave (1) with an added wave (β) equal to the reflected wave. This is not always true, but happens to be so for a scattering at one atom only. If there were a clump of impurity atoms, the amount added to the forward wave would not necessarily be the same as the reflected wave.

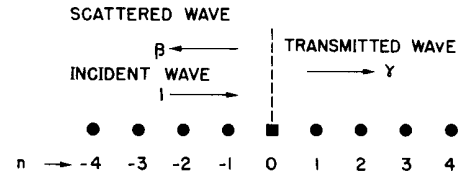


Fig. 13-6. Waves in a one-dimensional lattice with one “impurity” atom at $n = 0$.

We can get the amplitude β of the reflected wave from the middle equation of Eq. (13.33); we find that

$$\beta = \frac{-F}{F - 2iA \sin kb}. \quad (13.37)$$

We have the complete solution for the lattice with one unusual atom.

You may be wondering how the transmitted wave can be “more” than the incident wave as it appears in Eq. (13.34). Remember, though, that β and γ are complex numbers and that the number of particles (or rather, the probability of finding a particle) in a wave is proportional to the absolute square of the amplitude. In fact, there will be “conservation of electrons” only if

$$|\beta|^2 + |\gamma|^2 = 1. \quad (13.38)$$

You can show that this is true for our solution.

13-7 Trapping by a lattice imperfection

There is another interesting situation that can arise if F is a negative number. If the energy of the electron is lower at the impurity atom (at $n = 0$) than it is anywhere else, then the electron can get caught on this atom. That is, if $(E_0 + F)$ is below the bottom of the band at $(E_0 - 2A)$, then the electron can get “trapped” in a state with $E < E_0 - 2A$. Such a solution cannot come out of what we have done so far. We can get this solution, however, if we permit the trial solution we took in Eq. (13.15) to have an imaginary number for k . Let's set $k = i\kappa$. Again, we can have different solutions for $n < 0$ and for $n > 0$. A possible solution for $n < 0$ might be

$$a_n \text{ (for } n < 0) = ce^{+\kappa x_n}. \quad (13.39)$$

We have to take a plus sign in the exponent; otherwise the amplitude would get indefinitely large for large negative values of n . Similarly, a possible solution for $n > 0$ would be

$$a_n \text{ (for } n > 0) = c'e^{-\kappa x_n}. \quad (13.40)$$

If we put these trial solutions into Eq. (13.28) all but the middle three are satisfied provided that

$$E = E_0 - A(e^{\kappa b} + e^{-\kappa b}). \quad (13.41)$$

Since the sum of the two exponential terms is always greater than 2, this energy is below the regular band, and is what we are looking for. The remaining three equations in Eq. (13.28) are satisfied if $c = c'$ and if κ is chosen so that

$$A(e^{\kappa b} - e^{-\kappa b}) = -F. \quad (13.42)$$

Combining this equation with Eq. (13.41) we can find the energy of the trapped electron; we get

$$E = E_0 - \sqrt{4A^2 + F^2}. \quad (13.43)$$

The trapped electron has a unique energy—located somewhat below the conduction band.

Notice that the amplitudes we have in Eq. (13.39) and (13.40) do *not* say that the trapped electron sits right on the impurity atom. The probability of finding the electron at nearby atoms is given by the square of these amplitudes. For one particular choice of the parameters it might vary as shown in the bar graph of Fig. 13-7. The probability is greatest for finding the electron on the impurity atom. For nearby atoms the probability drops off exponentially with the distance from the impurity atom. This is another example of “barrier penetration.” From the point-of-view of classical physics the electron doesn't have enough energy to get away from the energy “hole” at the trapping center. But quantum mechanically it can leak out a little way.

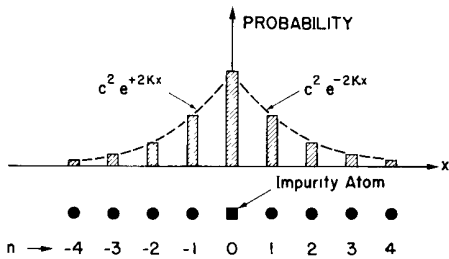


Fig. 13-7. The relative probabilities of finding a trapped electron at atomic sites near the trapping impurity atom.

13-8 Scattering amplitudes and bound states

Finally, our example can be used to illustrate a point which is very useful these days in the physics of high-energy particles. It has to do with a relationship between scattering amplitudes and bound states. Suppose we have discovered—through experiment and theoretical analysis—the way that pions scatter from protons. Then a new particle is discovered and someone wonders whether maybe it is just a combination of a pion and a proton held together in some bound state (in an analogy to the way an electron is bound to a proton to make a hydrogen atom). By a bound state we mean a combination which has a lower energy than the two free-particles.

There is a general theory which says that a bound state will exist at that energy at which the scattering amplitude becomes infinite if extrapolated algebraically (the mathematical term is “analytically continued”) to energy regions outside of the permitted band.

The physical reason for this is as follows. A bound state is a situation in which there are only waves tied on to a point and there’s no wave coming in to get it started, it just exists there by itself. The relative proportion between the so-called “scattered” or created wave and the wave being “sent in” is infinite. We can test this idea in our example. Let’s write our expression Eq. (13.37) for the scattered amplitude directly in terms of the energy E of the particle being scattered (instead of in terms of k). Since Equation (13.30) can be rewritten as

$$2A \sin kb = \sqrt{4A^2 - (E - E_0)^2},$$

the scattered amplitude is

$$\beta = \frac{-F}{F - i\sqrt{4A^2 - (E - E_0)^2}}. \quad (13.44)$$

From our derivation, this equation should be used only for real states—those with energies in the energy band, $E = E_0 \pm 2A$. But suppose we forget that fact and extend the formula into the “unphysical” energy regions where $|E - E_0| > 2A$. For these unphysical regions we can write†

$$\sqrt{4A^2 - (E - E_0)^2} = i\sqrt{(E - E_0)^2 - 4A^2}.$$

Then the “scattering amplitude,” whatever it may mean, is

$$\beta = \frac{-F}{F + \sqrt{(E - E_0)^2 - 4A^2}}. \quad (13.45)$$

Now we ask: Is there any energy E for which β becomes infinite (i.e., for which the expression for β has a “pole”)? Yes, so long as F is negative, the denominator of Eq. (13.45) will be zero when

$$(E - E_0)^2 - 4A^2 = F^2,$$

or when

$$E = E_0 \pm \sqrt{4A^2 + F^2}.$$

The minus sign gives just the energy we found in Eq. (13.43) for the trapped energy.

What about the plus sign? This gives an energy *above* the allowed energy band. And indeed there is another bound state there which we missed when we solved the equations of Eq. (13.28). We leave it as a puzzle for you to find the energy and amplitudes a_n for this bound state.

The relation between scattering and bound states provides one of the most useful clues in the current search for an understanding of the experimental observations about the new strange particles.

† The sign of the root to be chosen here is a technical point related to the allowed signs of κ in Eqs. (13.39) and (13.40). We won’t go into it here.

Semiconductors

14-1 Electrons and holes in semiconductors

One of the remarkable and dramatic developments in recent years has been the application of solid state science to technical developments in electrical devices such as transistors. The study of semiconductors led to the discovery of their useful properties and to a large number of practical applications. The field is changing so rapidly that what we tell you today may be incorrect next year. It will certainly be incomplete. And it is perfectly clear that with the continuing study of these materials many new and more wonderful things will be possible as time goes on. You will not need to understand this chapter for what comes later in this volume, but you may find it interesting to see that at least something of what you are learning has some relation to the practical world.

There are large numbers of semiconductors known, but we'll concentrate on those which now have the greatest technical application. They are also the ones that are best understood, and in understanding them we will obtain a degree of understanding of many of the others. The semiconductor substances in most common use today are silicon and germanium. These elements crystallize in the diamond lattice, a kind of cubic structure in which the atoms have tetrahedral bonding with their four nearest neighbors. They are insulators at very low temperatures—near absolute zero—although they do conduct electricity somewhat at room temperature. They are not metals; they are called *semiconductors*.

If we somehow put an extra electron into a crystal of silicon or germanium which is at a low temperature, we will have just the situation we described in the last chapter. The electron will be able to wander around in the crystal jumping from one atomic site to the next. Actually, we have looked only at the behavior of electrons in a rectangular lattice, and the equations would be somewhat different for the real lattice of silicon or germanium. All of the essential points are, however, illustrated by the results for the rectangular lattice.

As we saw in Chapter 13, these electrons can have energies only in a certain energy band—called the *conduction band*. Within this band the energy is related to the wave-number k of the probability amplitude C (see Eq. 13.24) by

$$E = E_0 - 2A_x \cos k_x a - 2A_y \cos k_y b - 2A_z \cos k_z c. \quad (14.1)$$

The A 's are the amplitudes for jumping in the x -, y -, and z -directions, and a , b , and c are the lattice spacings in these directions.

For energies near the bottom of the band, we can approximate Eq. (14.1) by

$$E = E_{\min} + A_x a^2 k_x^2 + A_y b^2 k_y^2 + A_z c^2 k_z^2 \quad (14.2)$$

(see Section 13-4)

If we think of electron motion in some particular direction, so that the components of k are always in the same ratio, the energy is a quadratic function of the wave number—and as we have seen of the momentum of the electron. We can write

$$E = E_{\min} + \alpha k^2, \quad (14.3)$$

where α is some constant, and we can make a graph of E versus k as in Fig. 14-1. We'll call such a graph an "energy diagram." An electron in a particular state of energy and momentum can be indicated by a point such as S in the figure

14-1 Electrons and holes in semiconductors

14-2 Impure semiconductors

14-3 The Hall effect

14-4 Semiconductor junctions

14-5 Rectification at a semiconductor junction

14-6 The transistor

Reference: C. Kittel, *Introduction to Solid State Physics*, Chapters 13, 14, and 18.

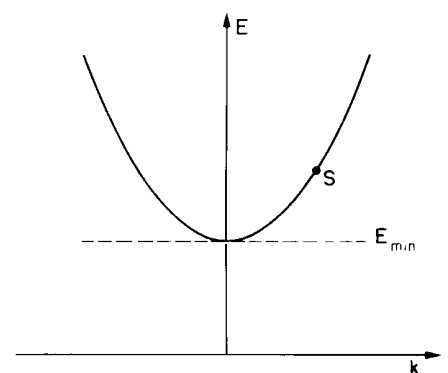


Fig. 14-1. The energy diagram for an electron in an insulating crystal.

As we also mentioned in Chapter 13, we can have a similar situation if we *remove* an electron from a neutral insulator. Then, an electron can jump over from a nearby atom and fill the “hole,” but leaving another “hole” at the atom it started from. We can describe this behavior by writing an amplitude to find the *hole* at any particular atom, and by saying that the *hole* can jump from one atom to the next. (Clearly, the amplitudes A that the hole jumps from atom a to atom b is just the same as the amplitude that an electron on atom b jumps into the hole at atom a .) The mathematics is just the same for the *hole* as it was for the extra electron, and we get again that the energy of the hole is related to its wave number by an equation just like Eq. (14.1) or (14.2), except, of course, with different numerical values for the amplitudes A_x , A_y , and A_z . The hole has an energy related to the wave number of its probability amplitudes. Its energy lies in a restricted band, and near the bottom of the band its energy varies quadratically with the wave number—or momentum—just as in Fig. 14-1. Following the arguments of Section 13-3, we would find that *the hole also behaves like a classical particle* with a certain effective mass—except that in noncubic crystals the mass depends on the direction of motion. So the hole behaves like a *positive particle* moving through the crystal. The charge of the hole-particle is positive, because it is located at the site of a missing electron; and when it moves in one direction there are actually electrons moving in the opposite direction.

If we put several electrons into a neutral crystal, they will move around much like the atoms of a low-pressure gas. If there are not too many, their interactions will not be very important. If we then put an electric field across the crystal, the electrons will start to move and an electric current will flow. Eventually they would all be drawn to one edge of the crystal, and, if there is a metal electrode there, they would be collected, leaving the crystal neutral.

Similarly we could put many holes into a crystal. They would roam around at random unless there is an electric field. With a field they would flow toward the negative terminal, and would be “collected”—what actually happens is that they are neutralized by electrons from the metal terminal.

One can also have both holes and electrons together. If there are not too many, they will all go their way independently. With an electric field, they will all contribute to the current. For obvious reasons, electrons are called the *negative carriers* and the holes are called the *positive carriers*.

We have so far considered that electrons are put into the crystal from the outside, or are removed to make a hole. It is also possible to “create” an electron-hole pair by taking a bound electron away from one neutral atom and putting it some distance away in the same crystal. We then have a free electron and a free hole, and the two can move about as we have described.

The energy required to put an electron *into* a state S —we say to “create” the state S —is the energy E^- shown in Fig. 14-2. It is some energy above E_{\min}^- . The energy required to “create” a hole in some state S' is the energy E^+ of Fig. 14-3, which is some energy greater than E_{\min}^+ . Now if we create a pair in the states S and S' , the energy required is just $E^- + E^+$.

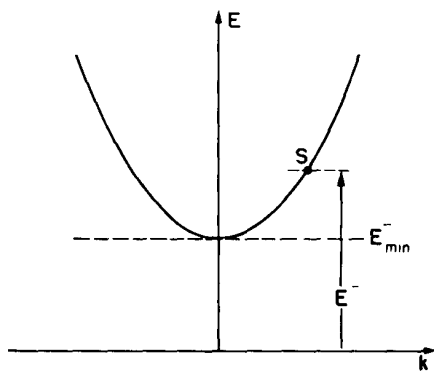


Fig. 14-2. The energy E^- is required to “create” a free electron.

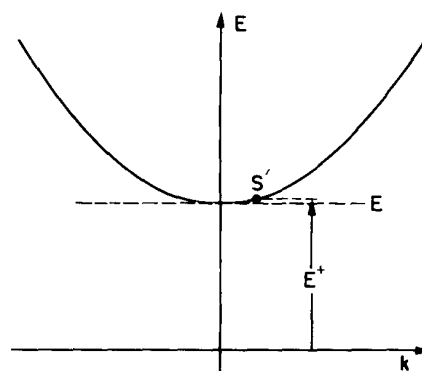


Fig. 14-3. The energy E^+ is required to “create” a hole in the state S' .

The creation of pairs is a common process (as we will see later), so many people like to put Fig. 14-2 and Fig. 14-3 together on the same graph—with the *hole* energy plotted *downward*, although it is, of course a *positive* energy. We have combined our two graphs in this way in Fig. 14-4. The advantage of such a graph is that the energy $E_{\text{pair}} = E^- + E^+$ required to create a pair with the electron in S and the hole in S' is just the vertical distance between S and S' as shown in Fig. 14-4. The minimum energy required to create a pair is called the “gap” energy and is equal to $E_{\text{min}}^- + E_{\text{min}}^+$.

Sometimes you will see a simpler diagram called an energy level diagram which is drawn when people are not interested in the k variable. Such a diagram—shown in Fig. 14-5—just shows the possible energies for the electrons and holes.†

How can electron-hole pairs be created? There are several ways. For example, photons of light (or x-rays) can be absorbed and create a pair if the photon energy is above the energy of the gap. The rate at which pairs are produced is proportional to the light intensity. If two electrodes are plated on a wafer of the crystal and a “bias” voltage is applied, the electrons and holes will be drawn to the electrodes. The circuit current will be proportional to the intensity of the light. This mechanism is responsible for the phenomenon of photoconductivity and the operation of photoconductive cells.

Electron hole pairs can also be produced by high-energy particles. When a fast-moving charged particle—for instance, a proton or a pion with an energy of tens or hundreds of Mev—goes through a crystal, its electric field will knock electrons out of their bound states creating electron-hole pairs. Such events occur hundreds of thousands of times per millimeter of track. After the passage of the particle, the carriers can be collected and in doing so will give an electrical pulse. This is the mechanism at play in the semiconductor counters recently put to use for experiments in nuclear physics. Such counters do not require semiconductors; they can also be made with crystalline insulators. In fact, the first of such counters was made using a diamond crystal which is an insulator at room temperature. Very pure crystals are required if the holes and electrons are to be able to move freely to the electrodes without being trapped. The semiconductors silicon and germanium are used because they can be produced with high purity in reasonable large sizes (centimeter dimensions).

So far we have been concerned with semiconductor crystals at temperatures near absolute zero. At any finite temperature there is still another mechanism by which electron-hole pairs can be created. The pair energy can be provided from the thermal energy of the crystal. The thermal vibrations of the crystal can transfer their energy to a pair—giving rise to “spontaneous” creation.

The probability per unit time that the energy as large as the gap energy E_{gap} will be concentrated at one atomic site is proportional to $e^{-E_{\text{gap}}/\kappa T}$, where T is the temperature and κ is Boltzmann’s constant (see Chapter 40, Vol. I). Near absolute zero there is no appreciable probability, but as the temperature rises there is an increasing probability of producing such pairs. At any finite temperature the production should continue forever at a constant rate giving more and more negative and positive carriers. Of course that does not happen because after awhile the electrons and holes accidentally find each other—the electron drops into the hole and the excess energy is given to the lattice. We say that the electron and hole “annihilate.” There is a certain probability per second that a hole meets an electron and the two things annihilate each other.

If the number of electrons per unit volume is N_n (n for negative carriers) and the density of positive carriers is N_p , the chance per unit time that an electron and a hole will find each other and annihilate is proportional to the product $N_n N_p$. In equilibrium this rate must equal the rate that pairs are created. You see that in

† In many books this same energy diagram is interpreted in a different way. The energy scale refers only to *electrons*. Instead of thinking of the energy of the hole, they think of the energy an electron *would* have if it filled the hole. This energy is *lower* than the free-electron energy—in fact, just the amount lower that you see in Fig. 14-5. With this interpretation of the energy scale, the gap energy is the minimum energy which must be given to an *electron* to move it from its bound state to the conduction band.

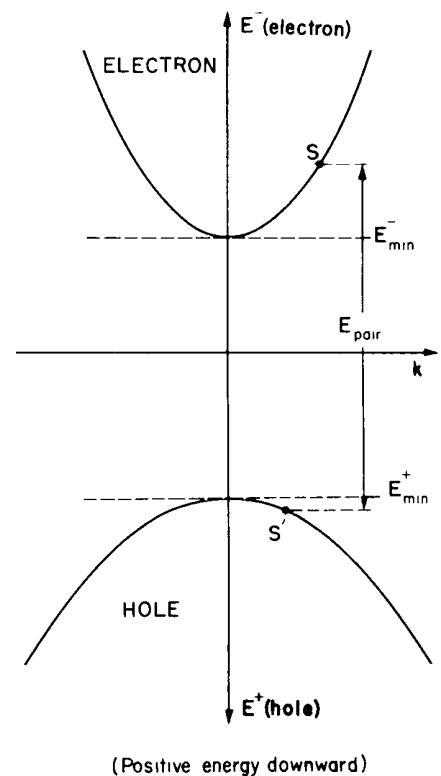


Fig. 14-4. Energy diagrams for an electron and a hole drawn together.

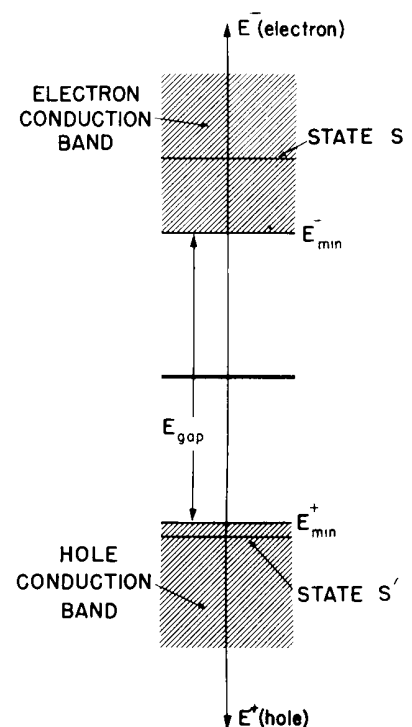


Fig. 14-5. Energy level diagram for electrons and holes.

equilibrium the product of N_n and N_p should be given by some constant times the Boltzmann factor:

$$N_n N_p = \text{const } e^{-E_{\text{gap}}/\kappa T}. \quad (14.4)$$

When we say constant, we mean nearly constant. A more complete theory—which includes more details about how holes and electrons “find” each other—shows that the “constant” is slightly dependent upon temperature, but the major dependence on temperature is in the exponential.

Let’s consider, as an example, a pure material which is originally neutral. At a finite temperature you would expect the number of positive and negative carriers to be equal, $N_n = N_p$. Then each of them should vary with temperature as $e^{-E_{\text{gap}}/2\kappa T}$. The variation of many of the properties of a semiconductor—the conductivity for example—is mainly determined by the exponential factor because all the other factors vary much more slowly with temperature. The gap energy for germanium is about 0.72 eV and for silicon 1.1 eV.

At room temperature κT is about 1/40 of an electron volt. At these temperatures there are enough holes and electrons to give a significant conductivity, while at, say, 30°K—one-tenth of room temperature—the conductivity is imperceptible. The gap energy of diamond is 6 or 7 eV and diamond is a good insulator at room temperature.

14-2 Impure semiconductors

So far we have talked about two ways that extra electrons can be put into an otherwise ideally perfect crystal lattice. One way was to inject the electron from an outside source; the other way, was to knock a bound electron off a neutral atom creating simultaneously an electron and a hole. It is possible to put electrons into the conduction band of a crystal in still another way. Suppose we imagine a crystal of germanium in which one of the germanium atoms is replaced by an arsenic atom. The germanium atoms have a valence of 4 and the crystal structure is controlled by the four valence electrons. Arsenic, on the other hand, has a valence of 5. It turns out that a single arsenic atom can sit in the germanium lattice (because it has approximately the correct size), but in doing so it must act as a valence 4 atom—using four of its valence electrons to form the crystal bonds and having one electron left over. This extra electron is very loosely attached—the binding energy is less than 1/10 of a volt. At room temperature the electron easily picks up that much energy from the thermal energy of the crystal, and then takes off on its own—moving about in the lattice as a free electron. An impurity atom such as the arsenic is called a *donor site* because it can give up a negative carrier to the crystal. If a crystal of germanium is grown from a melt to which a very small amount of arsenic has been added, the arsenic donor sites will be distributed throughout the crystal and the crystal will have a certain density of negative carriers built in.

You might think that these carriers would get swept away as soon as any small electric field was put across the crystal. This will not happen, however, because the arsenic atoms in the body of the crystal each have a positive charge. If the body of the crystal is to remain neutral, the average density of negative carrier electrons must be equal to the density of donor sites. If you put two electrodes on the edges of such a crystal and connect them to a battery, a current will flow; but as the carrier electrons are swept out at one end, new conduction electrons must be introduced from the electrode on the other end so that the average density of conduction electrons is left very nearly equal to the density of donor sites.

Since the donor sites are positively charged, there will be some tendency for them to capture some of the conduction electrons as they diffuse around inside the crystal. A donor site can, therefore, act as a trap such as those we discussed in the last section. But if the trapping energy is sufficiently small—as it is for arsenic—the number of carriers which are trapped at any one time is a small fraction of the total. For a complete understanding of the behavior of semiconductors

one must take into account this trapping. For the rest of our discussion, however, we will assume that the trapping energy is sufficiently low and the temperature is sufficiently high, that all of the donor sites have given up their electrons. This is, of course, just an approximation.

It is also possible to build into a germanium crystal some impurity atom whose valence is 3, such as aluminum. The aluminum atom tries to act as a valence 4 object by stealing an extra electron. It can steal an electron from some nearby germanium atom and end up as a negatively charged atom with an effective valence of 4. Of course, when it steals the electron from a germanium atom, it leaves a hole there; and this hole can wander around in the crystal as a positive carrier. An impurity atom which can produce a hole in this way is called an *acceptor* because it “accepts” an electron. If a germanium or a silicon crystal is grown from a melt to which a small amount of aluminum impurity has been added, the crystal will have built-in a certain density of holes which can act as positive carriers.

When a donor or an acceptor impurity is added to a semiconductor, we say that the material has been “doped.”

When a germanium crystal with some built-in donor impurities is at room temperature, some conduction electrons are contributed by the thermally induced electron-hole pair creation as well as by the donor sites. The electrons from both sources are, naturally, equivalent, and it is the total number N_n which comes into play in the statistical processes that lead to equilibrium. If the temperature is not too low, the number of negative carriers contributed by the donor impurity atoms is roughly equal to the number of impurity atoms present. In equilibrium Eq. (14.4) must still be valid; at a given temperature the product $N_n N_p$ is determined. This means that if we add some donor impurity which increases N_n , the number N_p of positive carriers will have to decrease by such an amount that $N_n N_p$ is unchanged. If the impurity concentration is high enough, the number N_n of negative carriers is determined by the number of donor sites and is nearly independent of temperature—all of the variation in the exponential factor is supplied by N_p , even though it is much less than N_n . An otherwise pure crystal with a small concentration of donor impurity will have a majority of negative carriers; such a material is called an “*n*-type” semiconductor.

If an acceptor-type impurity is added to the crystal lattice, some of the new holes will drift around and annihilate some of the free electrons produced by thermal fluctuation. This process will go on until Eq. (14.4) is satisfied. Under equilibrium conditions the number of positive carriers will be increased and the number of negative carriers will be decreased, leaving the product a constant. A material with an excess of positive carriers is called a “*p*-type” semiconductor.

If we put two electrodes on a piece of semiconductor crystal and connect them to a source of potential difference, there will be an electric field inside the crystal. The electric field will cause the positive and the negative carriers to move, and an electric current will flow. Let's consider first what will happen in an *n*-type material in which there is a large majority of negative carriers. For such material we can disregard the holes, they will contribute very little to the current because there are so few of them. In an ideal crystal the carriers would move across without any impediment. In a real crystal at a finite temperature, however,—especially in a crystal with some impurities—the electrons do not move completely freely. They are continually making collisions which knock them out of their original trajectories, that is, changing their momentum. These collisions are just exactly the scatterings we talked about in the last chapter and occur at any irregularity in the crystal lattice. In an *n*-type material the main causes of scattering are the very donor sites that are producing the carriers. Since the conduction electrons have a very slightly different energy at the donor sites, the probability waves are scattered from that point. Even in a perfectly pure crystal, however, there are (at any finite temperature) irregularities in the lattice due to thermal vibrations. From the classical point of view we can say that the atoms aren't lined up exactly on a regular lattice, but are, at any instant, slightly out of place due to their thermal

vibrations. The energy E_0 associated with each lattice point in the theory we described in Chapter 13 varies a little bit from place to place so that the waves of probability amplitude are not transmitted perfectly but are scattered in an irregular fashion. At very high temperatures or for very pure materials this scattering may become important, but in most doped materials used in practical devices the impurity atoms contribute most of the scattering. We would like now to make an estimate of the electrical conductivity of such a material.

When an electric field is applied to an n -type semiconductor, each negative carrier will be accelerated in this field, picking up velocity until it is scattered from one of the donor sites. This means that the carriers which are ordinarily moving about in a random fashion with their thermal energies will pick up an average drift velocity along the lines of the electric field and give rise to a current through the crystal. The drift velocity is in general rather small compared with the typical thermal velocities so that we can estimate the current by assuming that the average time that the carrier travels between scatterings is a constant. Let's say that the negative carrier has an effective electric charge q_n . In an electric field \mathbf{E} , the force on the carrier will be $q_n\mathbf{E}$. In Section 43-3 of Volume I we calculated the average drift velocity under such circumstances and found that it is given by $F\tau/m$, where F is the force on the charge, τ is the mean free time between collisions, and m is the mass. We should use the effective mass we calculated in the last chapter but since we want to make a rough calculation we will suppose that this effective mass is the same in all directions. Here we will call it m_n . With this approximation the average drift velocity will be

$$v_{\text{drift}} = \frac{q_n \mathbf{E} \tau_n}{m_n}. \quad (14.5)$$

Knowing the drift velocity we can find the current. Electric current density \mathbf{j} is just the number of carriers per unit volume, N_n , multiplied by the average drift velocity, and by the charge on each carrier. The current density is therefore

$$\mathbf{j} = N_n v_{\text{drift}} q_n \mathbf{E} = \frac{N_n q_n^2 \tau_n}{m_n} \mathbf{E}. \quad (14.6)$$

We see that the current density is proportional to the electric field; such a semiconductor material obeys Ohm's law. The coefficient of proportionality between \mathbf{j} and \mathbf{E} , the conductivity σ , is

$$\sigma = \frac{N_n q_n^2 \tau_n}{m_n}. \quad (14.7)$$

For an n -type material the conductivity is relatively independent of temperature. First, the number of majority carriers N_n is determined primarily by the density of donors in the crystal (so long as the temperature is not so low that too many of the carriers are trapped). Second, the mean time between collisions τ_n is mainly controlled by the density of impurity atoms, which is, of course, independent of the temperature.

We can apply all the same arguments to a p -type material, changing only the values of the parameters which appear in Eq. (14.7). If there are comparable numbers of both negative and positive carriers present at the same time, we must add the contributions from each kind of carrier. The total conductivity will be given by

$$\sigma = \frac{N_n q_n^2 \tau_n}{m_n} + \frac{N_p q_p^2 \tau_p}{m_p}. \quad (14.8)$$

For very pure materials, N_p and N_n will be nearly equal. They will be smaller than in a doped material, so the conductivity will be less. Also they will vary rapidly with temperature (like $e^{-E_{\text{gap}}/\kappa T}$, as we have seen), so the conductivity may change extremely fast with temperature.

14-3 The Hall effect

It is certainly a peculiar thing that in a substance where the only relatively free objects are electrons, there should be an electrical current carried by holes that behave like positive particles. We would like, therefore, to describe an experiment that shows in a rather clear way that the sign of the carrier of electric current is quite definitely positive. Suppose we have a block made of semiconductor material—it could also be a metal—and we put an electric field on it so as to draw a current in some direction, say the horizontal direction as drawn in Fig. 14-6. Now suppose we put a magnetic field on the block pointing at a right angle to the current, say *into* the plane of the figure. The moving carriers will feel a magnetic force $q(\mathbf{v} \times \mathbf{B})$. And since the average drift velocity is either right or left—depending on the sign of the charge on the carrier—the average magnetic force on the carriers will be either up or down. No, that is not right! For the directions we have assumed for the current and the magnetic field the magnetic force on the moving charges will always be *up*. Positive charges moving in the direction of \mathbf{j} (to the right) will feel an upward force. If the current is carried by negative charges, they will be moving left (for the same sign of the conduction current) and they will also feel an upward force. Under steady conditions, however, there is no upward motion of the carriers because the current can flow only from left to right. What happens is that a few of the charges initially flow upward, producing a surface charge density along the upper surface of semiconductor—leaving an equal and opposite surface charge density along the bottom surface of the crystal. The charges pile up on the top and bottom surfaces until the electric forces they produce on the moving charges just exactly cancel the magnetic force (on the average) so that the steady current flows horizontally. The charges on the top and bottom surfaces will produce a potential difference vertically across the crystal which can be measured with a high-resistance voltmeter, as shown in Fig. 14-7. The sign of the potential difference registered by the voltmeter will depend on the sign of the carrier charges responsible for the current.

When such experiments were first done it was expected that the sign of the potential difference would be negative as one would expect for negative conduction electrons. People were, therefore, quite surprised to find that for some materials the sign of the potential difference was in the opposite direction. It appeared that the current carrier was a particle with a positive charge. From our discussion of doped semiconductors it is understandable that an *n*-type semiconductor should produce the sign of potential difference appropriate to negative carriers, and that a *p*-type semiconductor should give an opposite potential difference, since the current is carried by the positively charged holes.

The original discovery of the anomalous sign of the potential difference in the Hall effect was made in a metal rather than a semiconductor. It had been assumed that in metals the conduction was always by electron; however, it was found out that for beryllium the potential difference had the wrong sign. It is now understood that in metals as well as in semiconductors it is possible, in certain circumstances, that the “objects” responsible for the conduction are holes. Although it is ultimately the electrons in the crystal which do the moving, nevertheless, the relationship of the momentum and the energy, and the response to external fields is exactly what one would expect for an electric current carried by positive particles.

Let’s see if we can make a quantitative estimate of the magnitude of the voltage difference expected from the Hall effect. If the voltmeter in Fig 14-7 draws a negligible current, then the charges inside the semiconductor must be moving from left to right and the vertical magnetic force must be precisely cancelled by a vertical electric field which we will call \mathcal{E}_{tr} (the “tr” is for “transverse”). If this electric field is to cancel the magnetic forces, we must have

$$\mathcal{E}_{tr} = -v_{drift} \times B. \quad (14.9)$$

Using the relation between the drift velocity and the electric current density given

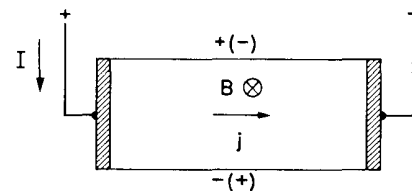


Fig. 14-6. The Hall effect comes from the magnetic forces on the carriers.

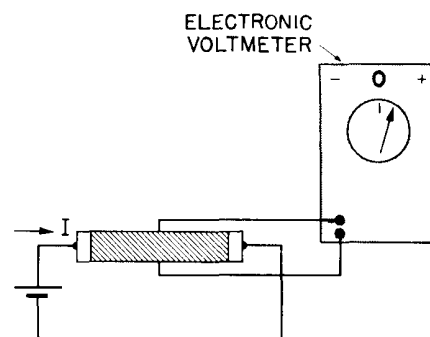


Fig. 14-7. Measuring the Hall effect.

in Eq. (14.6), we get

$$\epsilon_{tr} = -\frac{1}{qN} jB.$$

The potential difference between the top and the bottom of the crystal is, of course, this electric field strength multiplied by the height of the crystal. The electric field strength ϵ_{tr} in the crystal is proportional to the current density and to the magnetic field strength. The constant of proportionality $1/qN$ is called the Hall coefficient and is usually represented by the symbol R_H . The Hall coefficient depends just on the density of carriers—provided that carriers of one sign are in a large majority. Measurement of the Hall effect is, therefore, one convenient way of determining experimentally the density of carriers in a semiconductor.

14-4 Semiconductor junctions

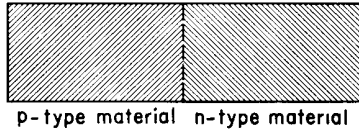


Fig. 14-8. A p - n junction

We would like to discuss now what happens if we take two pieces of germanium or silicon with different internal characteristics—say different kinds or amounts of doping—and put them together to make a “junction.” Let’s start out with what is called a p - n junction in which we have p -type germanium on one side of the boundary and n -type germanium on the other side of the boundary—as sketched in Fig. 14-8. Actually, it is not practical to put together two separate pieces of crystal and have them in uniform contact on an atomic scale. Instead, junctions are made out of a single crystal which has been modified in the two separate regions. One way is to add some suitable doping impurity to the “melt” after only half of the crystal has grown. Another way is to paint a little of the impurity element on the surface and then heat the crystal causing some impurity atoms to diffuse into the body of the crystal. Junctions made in these ways do not have a sharp boundary, although the boundaries can be made as thin as 10^{-4} centimeters or so. For our discussions we will imagine an ideal situation in which these two regions of the crystal with different properties meeting at a sharp boundary.

On the n -type side of p - n junction there are free electrons which can move about, as well as the fixed donor sites which balance the overall electric charge. On the p -type side there are free holes moving about and an equal number of negative acceptor sites keeping the charge balanced. Actually, that describes the situation before we put the two materials in contact. Once they are connected together the situation will change near the boundary. When the electrons in the n -type material arrive at the boundary they will not be reflected back as they would at a free surface, but are able to go right on into the p -type material. Some of the electrons of the n -type material will, therefore, tend to diffuse over into the p -type material where there are fewer electrons. This cannot go on forever because as we lose electrons from the n -side the net positive charge there increases until finally an electric voltage is built up which retards the diffusion of electrons into the p -side. In a similar way, the positive carriers of the p -type material can diffuse across the junction into the n -type material. When they do this they leave behind an excess of negative charge. Under equilibrium conditions the net diffusion current must be zero. This brought about by the electric fields which are established in such a way as to draw the positive carriers back toward the p -type material.

The two diffusion processes we have been describing go on simultaneously and, you will notice, both act in the direction which will charge up the n -type material in a positive sense and the p -type material in a negative sense. Because of the finite conductivity of the semiconductor material, the change in potential from the p -side to the n -side will occur in a relatively narrow region near the boundary; the main body of each block of material will have a uniform potential. Let’s imagine an x -axis in a direction perpendicular to the boundary surface. Then the electric potential will vary with x , as shown in Fig. 14-9(b). We have also shown in part (c) of the figure the expected variation of the density N_n of n -carriers and the density N_p of p -carriers. Far away from the junction the carrier densities N_p and N_n should be just the equilibrium density we would expect for individual blocks of materials at the same temperature. (We have drawn the figure for a

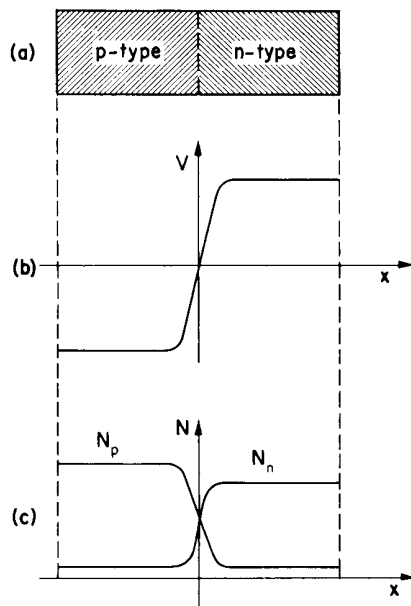


Fig. 14-9. The electric potential and the carrier densities in an unbiased semiconductor junction.

junction in which the p -type material is more heavily doped than the n -type material.) Because of the potential gradient at the junction, the positive carriers have to climb up a potential hill to get to the p -type side. This means that under equilibrium conditions there can be fewer positive carriers in the n -type material than there are in the p -type material. Remembering the laws of statistical mechanics, we expect that the ratio of p -type carriers on the two sides to be given by the following equation:

$$\frac{N_p(p\text{-side})}{N_p(n\text{-side})} = e^{-q_p V / \kappa T}. \quad (14.10)$$

The product $q_p V$ in the numerator of the exponential is just the energy required to carry a charge of q_p through a potential difference V .

We have a precisely similar equation for the densities of the n -type carriers:

$$\frac{N_n(n\text{-side})}{N_n(p\text{-side})} = e^{-q_n V / \kappa T}. \quad (14.11)$$

If we know the equilibrium densities in each of the two materials, we can use either of the two equations above to determine the potential difference across the junction.

Notice that if Eqs. (14.10) and (14.11) are to give the same value for the potential difference V , the product $N_p N_n$ must be the same for the p -side as for the n -side. (Remember that $q_n = -q_p$.) We have seen earlier, however, that this product depends only on the temperature and the gap energy of the crystal. Provided both sides of the crystal are at the same temperature, the two equations are consistent with the same value of the potential difference.

Since there is a potential difference from one side of the junction to the other, it looks something like a battery. Perhaps if we connect a wire from the n -type side to the p -type side we will get an electrical current. That would be nice because then the current would flow forever without using up any material and we would have an infinite source of energy in violation of the second law of thermodynamics! There is, however, no current if you connect a wire from the p -side to the n -side. And the reason is easy to see. Suppose we imagine first a wire made out of a piece of undoped material. When we connect this wire to the n -type side, we have a junction. There will be a potential difference across this junction. Let's say that it is just one-half the potential difference from the p -type material to the n -type material. When we connect our undoped wire to the p -type side of the junction, there is also a potential difference at this junction—again, one-half the potential drop across the p - n junction. At all the junctions the potential differences adjust themselves so that there is no net current flow in the circuit. Whatever kind of wire you use to connect together the two sides of the n - p junction, you are producing two new junctions, and so long as all the junctions are at the same temperature, the potential jumps at the junctions all compensate each other and no current will flow in the circuit. It does turn out, however—if you work out the details—that if some of the junctions are at a different temperature than the other junctions, currents will flow. Some of the junctions will be heated and others will be cooled by this current and thermal energy will be converted into electrical energy. This effect is responsible for the operation of thermocouples which are used for measuring temperatures, and of thermoelectric generators. The same effect is also used to make small refrigerators.

If we cannot measure the potential difference between the two sides of an n - p junction, how can we really be sure that the potential gradient shown in Fig. 14-9 really exists? One way is to shine light on the junction. When the light photons are absorbed they can produce an electron-hole pair. In the strong electric field that exists at the junction (equal to the slope of the potential curve of Fig. 14-9) the hole will be driven into the p -type region and the electron will be driven into the n -type region. If the two sides of the junction are now connected to an external circuit, these extra charges will provide a current. The energy of the light will be converted into electrical energy in the junction. The solar cells which generate electrical power for the operation of some of our satellites operate on this principle.

In our discussion of the operation of a semiconductor junction we have been assuming that the holes and the electrons act more-or-less independently—except that they somehow get into proper statistical equilibrium. When we were describing the current produced by light shining on the junction, we were assuming that an electron or a hole produced in the junction region would get into the main body of the crystal before being annihilated by a carrier of the opposite polarity. In the immediate vicinity of the junction, where the density of carriers of both signs is approximately equal, the effect of electron-hole annihilation (or as it is often called, “recombination”) is an important effect, and in a detailed analysis of a semiconductor junction must be properly taken into account. We have been assuming that a hole or an electron produced in a junction region has a good chance of getting into the main body of the crystal before recombining. The typical time for an electron or a hole to find an opposite partner and annihilate it is for typical semiconductor materials in the range between 10^{-3} and 10^{-7} seconds. This time is, incidentally, much longer than the mean free time τ between collisions with scattering sites in the crystal which we used in the analysis of conductivity. In a typical n - p junction, the time for an electron or hole formed in the junction region to be swept away into the body of the crystal is generally much shorter than the recombination time. Most of the pairs will, therefore, contribute to an external current.

14-5 Rectification at a semiconductor junction

We would like to show next how it is that a p - n junction can act like a rectifier. If we put a voltage across the junction, a large current will flow if the polarity is in one direction, but a very small current will flow if the same voltage is applied in the opposite direction. If an alternating voltage is applied across the junction, a net current will flow in one direction—the current is “rectified.” Let’s look again at what is going on in the equilibrium condition described by the graphs of Fig. 14-9. In the p -type material there is a large concentration N_p of positive carriers. These carriers are diffusing around and a certain number of them each second approach the junction. This current of positive carriers which approaches the junction is proportional to N_p . Most of them, however, are turned back by the high potential hill at the junction and only the fraction $e^{-qV/\kappa T}$ gets through. There is also a current of positive carriers approaching the junction from the other side. This current is also proportional to the density of positive carriers in the n -type region, but the carrier density here is much smaller than the density on the p -type side. When the positive carriers approach the junction from the n -type side, they find a hill with a negative slope and immediately slide downhill to the p -type side of the junction. Let’s call this current I_0 . Under equilibrium the currents from the two directions are equal. We expect then the following relation:

$$I_0 \sim N_p(n\text{-side}) = N_p(p\text{-side})e^{-qV/\kappa T}. \quad (14.12)$$

You will notice that this equation is really just the same as Eq. (14-10). We have just derived it in a different way.

Suppose, however, that we lower the voltage on the n -side of the junction by an amount ΔV —which we can do by applying an external potential difference to the junction. Now the difference in potential across the potential hill is no longer V but $V - \Delta V$. The current of positive carriers from the p -side to the n -side will now have this potential difference in its exponential factor. Calling this current I_1 , we have

$$I_1 \sim N_p(p\text{-side})e^{-q(V-\Delta V)/\kappa T}.$$

This current is larger than I_0 by just the factor $e^{q\Delta V/\kappa T}$. So we have the following relation between I_1 and I_0 :

$$I_1 = I_0 e^{q\Delta V/\kappa T}. \quad (14.13)$$

The current from the p -side increases exponentially with the externally applied voltage ΔV . The current of positive carriers from the n -side, however, remains

constant so long as ΔV is not too large. When they approach the barrier, these carriers will still find a downhill potential and will all fall down to the p -side (If ΔV is larger than the natural potential difference V , the situation would change, but we will not consider what happens at such high voltages.) The net current I of positive carriers which flows across the junction is then the difference between the currents from the two sides:

$$I = I_0(e^{+q\Delta V/kT} - 1). \quad (14.14)$$

The net current I of holes flows into the n -type region. There the holes diffuse into the body of the n -region, where they are eventually annihilated by the majority n -type carriers—the electrons. The electrons which are lost in this annihilation will be made up by a current of electrons from the external terminal of the n -type material.

When ΔV is zero, the net current in Eq. (14.14) is zero. For positive ΔV the current increases rapidly with the applied voltage. For negative ΔV the current reverses in sign, but the exponential term soon becomes negligible and the negative current never exceeds I_0 —which under our assumptions is rather small. This back current I_0 is limited by the small density of the minority carriers on the n -side of the junction.

If you go through exactly the same analysis for the current of negative carriers which flows across the junction, first with no potential difference and then with a small externally applied potential difference ΔV , you get again an equation just like (14.14) for the net electron current. Since the total current is the sum of the currents contributed by the two carriers, Eq. (14.14) still applies for the total current provided we identify I_0 as the maximum current which can flow for a reversed voltage.

The voltage-current characteristic of Eq. (14.14) is shown in Fig. 14-10. It shows the typical behavior of solid state diodes—such as those used in modern computers. We should remark that Eq. (14.14) is true only for small voltages. For voltages comparable to or larger than the natural internal voltage difference V , other effects come into play and the current no longer obeys the simple equation.

You may remember, incidentally, that we got exactly the same equation we have found here in Eq. (14.14) when we discussed the “mechanical rectifier”—the ratchet and pawl—in Chapter 46 of Volume I. We get the same equations in the two situations because the basic physical processes are quite similar.

14-6 The transistor

Perhaps the most important application of semiconductors is in the transistor. The transistor consists of two semiconductor junctions very close together. Its operation is based in part on the same principles that we just described for the semiconductor diode—the rectifying junction. Suppose we make a little bar of germanium with three distinct regions, a p -type region, an n -type region, and another p -type region, as shown in Fig. 14-11(a). This combination is called a p - n - p transistor. Each of the two junctions in the transistor will behave much in the way we have described in the last section. In particular, there will be a potential gradient at each junction having a certain potential drop from the n -type region to each p -type region. If the two p -type regions have the same internal properties, the variation in potential as we go across the crystal will be as shown in the graph of Fig. 14-11(b).

Now let's imagine that we connect each of the three regions to external voltage sources as shown in part (a) of Fig. 14-12. We will refer all voltages to the terminal connected to the left-hand p -region so it will be, by definition, at zero potential. We will call this terminal the *emitter*. The n -type region is called the *base* and it is connected to a slightly negative potential. The right-hand p -type region is called the *collector*, and is connected to a somewhat larger negative potential. Under these circumstances the variation of potential across the crystal will be as shown in the graph of Fig. 14-12(b).

Let's first see what happens to the positive carriers, since it is primarily their behavior which controls the operation of the p - n - p transistor. Since the emitter is

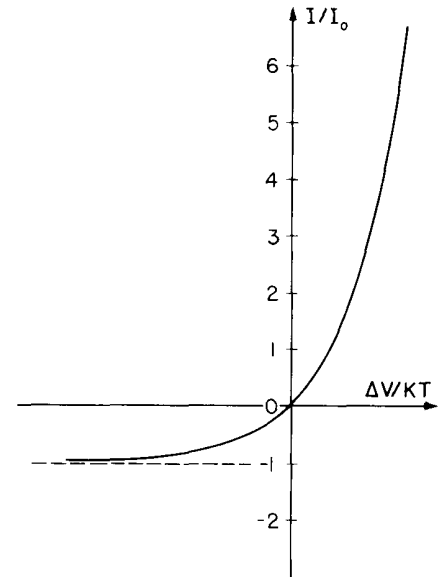


Fig. 14-10. The current through a junction as a function of the voltage across it.

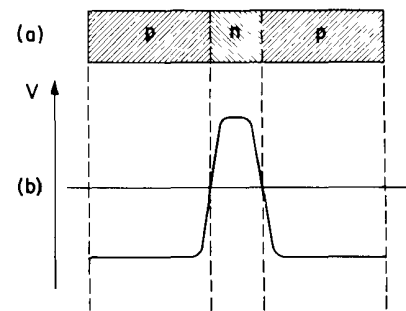


Fig. 14-11. The potential distribution in a transistor with no applied voltages.

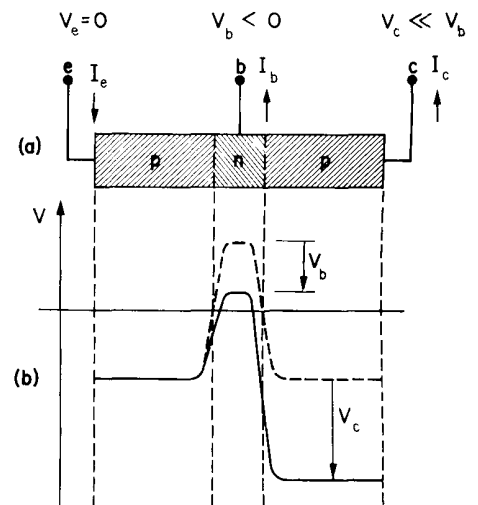


Fig. 14-12. The potential distribution in an operating transistor.

at a relatively more positive potential than the base, a current of positive carriers will flow from the emitter region into the base region. A relatively large current flows, since we have a junction operating with a “forward voltage”—corresponding to the right-hand half of the graph in Fig. 14-10. With these conditions, positive carriers or holes are being “emitted” from the p -type region into the n -type region. You might think that this current would flow out of the n -type region through the base terminal b . Now, however, comes the secret of the transistor. The n -type region is made very thin—typically 10^{-3} cm or less, much narrower than its transverse dimensions. This means that as the holes enter the n -type region they have a very good chance of diffusing across to the other junction before they are annihilated by the electrons in the n -type region. When they get to the right-hand boundary of the n -type region they find a steep downward potential hill and immediately fall into the right-hand p -type region. This side of the crystal is called the collector because it “collects” the holes after they have diffused across the n -type region. In a typical transistor, all but a fraction of a percent of the hole current which leaves the emitter and enters the base is collected in the collector region, and only the small remainder contributes to the net base current. The sum of the base and collector currents is, of course, equal to the emitter current.

Now imagine what happens if we vary slightly the potential V_b on the base terminal. Since we are on a relatively steep part of the curve of Fig. 14-10, a small variation of the potential V_b will cause a rather large change in the emitter current I_e . Since the collector voltage V_c is much more negative than the base voltage, these slight variations in potential will not effect appreciably the steep potential hill between the base and the collector. Most of the positive carriers emitted into the n -region will still be caught by the collector. Thus as we vary the potential of the base electrode, there will be a corresponding variation in the collector current I_c . The essential point, however, is that the base current I_b always remains a small fraction of the collector current. The transistor is an amplifier; a small current I_b introduced into the base electrode gives a large current—100 or so times higher—at the collector electrode.

What about the electrons—the negative carriers that we have been neglecting so far? First, note that we do not expect any significant electron current to flow between the base and the collector. With a large negative voltage on the collector, the electrons in the base would have to climb a very high potential energy hill and the probability of doing that is very small. There is a very small current of electrons to the collector.

On the other hand, the electrons in the base *can* go into the emitter region. In fact, you might expect the electron current in this direction to be comparable to the hole current from the emitter into the base. Such an electron current isn't useful, and, on the contrary, is bad because it increases the total base current required for a given current of holes to the collector. The transistor is, therefore, designed to minimize the electron current to the emitter. The electron current is proportional to $N_n(\text{base})$, the density of negative carriers in the base material while the hole current from the emitter depends on $N_p(\text{emitter})$, the density of positive carriers in the emitter region. By using relatively little doping in the n -type material $N_n(\text{base})$ can be made much smaller than $N_p(\text{emitter})$. (The very thin base region also helps a great deal because the sweeping out of the holes in this region by the collector increases significantly the average hole current from the emitter into the base, while leaving the electron current unchanged.) The net result is that the electron current across the emitter-base junction can be made much less than the hole current, so that the electrons do not play any significant role in operation of the p - n - p transistor. The currents are dominated by motion of the holes, and the transistor performs as an amplifier as we have described above.

It is also possible to make a transistor by interchanging the p -type and n -type materials in Fig. 14-11. Then we have what is called an n - p - n transistor. In the n - p - n transistor the main currents are carried by the electrons which flow from the emitter into the base and from there to the collector. Obviously, all the arguments we have made for the p - n - p transistor also apply to the n - p - n transistor if the potentials of the electrodes are chosen with the opposite signs.

The Independent Particle Approximation

15-1 Spin waves

In Chapter 13 we worked out the theory for the propagation of an electron or of some other “particle,” such as an atomic excitation, through a crystal lattice. In the last chapter we applied the theory to semiconductors. But when we talked about situations in which there are many electrons we disregarded any interactions between them. To do this is of course only an approximation. In this chapter we will discuss further the idea that you can disregard the interaction between the electrons. We will also use the opportunity to show you some more applications of the theory of the propagation of particles. Since we will generally continue to disregard the interactions between particles, there is very little really new in this chapter except for the new applications. The first example to be considered is, however, one in which it is possible to write down quite exactly the correct equations when there is more than one “particle” present. From them we will be able to see how the approximation of disregarding the interactions is made. We will not, though, analyze the problem very carefully.

As our first example we will consider a “spin wave” in a ferromagnetic crystal. We have discussed the theory of ferromagnetism in Chapter 36 of Volume II. At zero temperature all the electron spins that contribute to the magnetism in the body of a ferromagnetic crystal are parallel. There is an interaction energy between the spins, which is lowest when all the spins are down. At any nonzero temperature, however, there is some chance that some of the spins are turned over. We calculated the probability in an approximate manner in Chapter 36. This time we will describe the quantum mechanical theory—so you will see what you would have to do if you wanted to solve the problem more exactly. (We will still make some idealizations by assuming that the electrons are localized at the atoms and that the spins interact only with neighboring spins.)

We consider a model in which the electrons at each atom are all paired except one, so that all of the magnetic effects come from one spin- $\frac{1}{2}$ electron per atom. Further, we imagine that these electrons are localized at the atomic sites in the lattice. The model corresponds roughly to metallic nickel.

We also assume that there is an interaction between any two adjacent spinning electrons which gives a term in the energy of the system

$$E = - \sum_{i,j} K \sigma_i \cdot \sigma_j, \quad (15.1)$$

where σ 's represent the spins and the summation is over all adjacent pairs of electrons. We have already discussed this kind of interaction energy when we considered the hyperfine splitting of hydrogen due to the interaction of the magnetic moments of the electron and proton in a hydrogen atom. We expressed it then as $A \sigma_e \cdot \sigma_p$. Now, for a given pair, say the electrons at atom 4 and at atom 5, the Hamiltonian would be $-K \sigma_4 \cdot \sigma_5$. We have a term for each such pair, and the Hamiltonian is (as you would expect for classical energies) the sum of these terms for each interacting pair. The energy is written with the factor $-K$ so that a positive K will correspond to ferromagnetism—that is, the lowest energy results when adjacent spins are parallel. In a real crystal, there may be other terms which are the interactions of *next* nearest neighbors, and so on, but we don't need to consider such complications at this stage.

With the Hamiltonian of Eq. (15.1) we have a complete description of the ferromagnet—within our approximation—and the properties of the magnetization

15-1 Spin waves

15-2 Two spin waves

15-3 Independent particles

15-4 The benzene molecule

15-5 More organic chemistry

15-6 Other uses of the approximation

should come out. We should also be able to calculate the thermodynamic properties due to the magnetization. If we can find all the energy levels, the properties of the crystal at a temperature T can be found from the principle that the probability that a system will be found in a given state of energy E is proportional to $e^{-E/kT}$. This problem has never been completely solved.

We will show some of the problems by taking a simple example in which all the atoms are in a line—a one-dimensional lattice. You can easily extend the ideas to three dimensions. At each atomic location there is an electron which has two possible states, either spin up or spin down, and the whole system is described by telling how all of the spins are arranged. We take the Hamiltonian of the system to be the operator of the interaction energy. Interpreting the spin vectors of Eq. (15.1) as the sigma-operators—or the sigma-matrices—we write for the linear lattice

$$\hat{H} = \sum_n -\frac{A}{2} \hat{\sigma}_n \cdot \hat{\sigma}_{n+1}. \quad (15.2)$$

In this equation we have written the constant as $A/2$ for convenience (so that some of the later equations will be exactly the same as the ones in Chapter 13).

Now what is the lowest state of this system? The state of lowest energy is the one in which all the spins are parallel—let's say, all up.† We can write this state as $|\cdots + + + \cdots\rangle$, or $|\text{gnd}\rangle$ for the “ground,” or lowest, state. It's easy to figure out the energy for this state. One way is to write out all the vector sigmas in terms of $\hat{\sigma}_x$, $\hat{\sigma}_y$, and $\hat{\sigma}_z$, and work through carefully what each term of the Hamiltonian does to the ground state, and then add the results. We can, however, also use a good short cut. We saw in Section 12-2, that $\hat{\sigma}_i \cdot \hat{\sigma}_j$ could be written in terms of the Pauli spin exchange operator like this:

$$\hat{\sigma}_i \cdot \hat{\sigma}_j = (2\hat{P}_{ij}^{\text{spin ex}} - 1), \quad (15.3)$$

where the operator $\hat{P}_{ij}^{\text{spin ex}}$ interchanges the spins of the i th and j th electrons. With this substitution the Hamiltonian becomes

$$\hat{H} = -A \sum_n (\hat{P}_{n,n+1}^{\text{spin ex}} - \frac{1}{2}). \quad (15.4)$$

It is now easy to work out what happens to different states. For instance if i and j are both up, then exchanging the spins leaves everything unchanged, so \hat{P}_{ij} acting on the state just gives the same state back, and is equivalent to multiplying by $+1$. The expression $(\hat{P}_{ij} - \frac{1}{2})$ is just equal to one-half. (From now on we will leave off the descriptive superscript on the P .)

For the ground state all spins are up; so if you exchange a particular pair of spins, you get back the original state. The ground state is a stationary state. If you operate on it with the Hamiltonian you get the same state again multiplied by a sum of terms, $-(A/2)$ for each pair of spins. That is, the energy of the system in the ground state is $-A/2$ per atom

Next we would like to look at the energies of some of the excited states. It will be convenient to measure the energies with respect to the ground state—that is, to choose the ground state as our zero of energy. We can do that by adding the energy $A/2$ to each term in the Hamiltonian. That just changes the “ $\frac{1}{2}$ ” in Eq. (15.4) to “1.” Our new Hamiltonian is

$$\hat{H} = -A \sum_n (\hat{P}_{n,n+1} - 1). \quad (15.5)$$

With this Hamiltonian the energy of the lowest state is zero; the spin exchange operator is equivalent to multiplying by unity (for the ground state) which is cancelled by the “1” in each term.

† The ground state here is really “degenerate”; there are other states with the same energy—for example, all spins down, or all in any other direction. The slightest external field in the z -direction will give a different energy to all these states, and the one we have chosen will be the true ground state.

For describing states other than the ground state we will need a suitable set of base states. One convenient approach is to group the states according to whether one electron has spin down, or two electrons have spin down, and so on. There are, of course, many states with one spin down. The down spin could be at atom "4," or at atom "5," or at atom "6," . . . We can, in fact, choose just such states for our base states. We could write them this way: $|4\rangle$, $|5\rangle$, $|6\rangle$, . . . It will, however, be more convenient later if we label the "odd atom"—the one with the down-spinning electron—by its coordinate x . That is, we'll define the state $|x_5\rangle$ to be one with all the electrons spinning up except for the one on the atom at x_5 , which has a down-spinning electron (see Fig. 15-1). In general, $|x_n\rangle$ is the state with one down spin that is located at the coordinate x_n of the n th atom.

What is the action of the Hamiltonian (15.5) on the state $|x_5\rangle$? One term of the Hamiltonian is say $-A(\hat{P}_{7,8} - 1)$. The operator $\hat{P}_{7,8}$ exchanges the two spins of the adjacent atoms 7, 8. But in the state $|x_5\rangle$ these are both up, and nothing happens; $\hat{P}_{7,8}$ is equivalent to multiplying by 1:

$$\hat{P}_{7,8} |x_5\rangle = |x_5\rangle.$$

It follows that

$$(\hat{P}_{7,8} - 1) |x_5\rangle = 0.$$

Thus all the terms of the Hamiltonian give zero—except those involving atom 5, of course. On the state $|5\rangle$, the operation $\hat{P}_{4,5}$ exchanges the spin of atom 4 (up) and atom 5 (down). The result is the state with all spins up except the atom at 4. That is

$$\hat{P}_{4,5} |x_5\rangle = |x_4\rangle.$$

In the same way

$$\hat{P}_{5,6} |x_5\rangle = |x_6\rangle.$$

Hence, the only terms of the Hamiltonian which survive are $-A(\hat{P}_{4,5} - 1)$ and $-A(\hat{P}_{5,6} - 1)$. Acting on $|x_5\rangle$ they produce $-A|x_4\rangle + A|x_5\rangle$ and $-A|x_6\rangle + A|x_5\rangle$, respectively. The result is

$$\hat{H} |x_5\rangle = -A \sum_n (\hat{P}_{n,n+1} - 1) |x_5\rangle = -A\{|x_6\rangle + |x_4\rangle - 2|x_5\rangle\}. \quad (15.6)$$

When the Hamiltonian acts on state $|x_5\rangle$ it gives rise to some amplitude to be in states $|x_4\rangle$ and $|x_6\rangle$. That just means that there is a certain amplitude to have the down spin jump over to the next atom. So because of the interaction between the spins, if we begin with one spin down, then there is some probability that at a later time another one will be down instead. Operating on the general state $|x_n\rangle$, the Hamiltonian gives

$$\hat{H} |x_n\rangle = -A\{|x_{n+1}\rangle + |x_{n-1}\rangle - 2|x_n\rangle\}. \quad (15.7)$$

Notice particularly that if we take a complete set of states with only one spin down, they will only be mixed among themselves. The Hamiltonian will never mix these states with others that have more spins down. So long as you only exchange spins you never change the total number of down spins.

It will be convenient to use the matrix notation for the Hamiltonian, say $H_{n,m} \equiv \langle x_n | \hat{H} | x_m \rangle$; Eq. (15.7) is equivalent to

$$\begin{aligned} H_{n,n} &= A; \\ H_{n,n+1} &= H_{n,n-1} = -A; \\ H_{n,m} &= 0, \quad \text{for } |n - m| > 1. \end{aligned} \quad (15.8)$$

Now what are the energy levels for states with one spin down? As usual we let C_n be the amplitude that some state $|\psi\rangle$ is in the state $|x_n\rangle$. If $|\psi\rangle$ is to be a definite energy state, all the C_n 's must vary with time in the same way, namely,

$$C_n = a_n e^{-iEt/\hbar}. \quad (15.9)$$

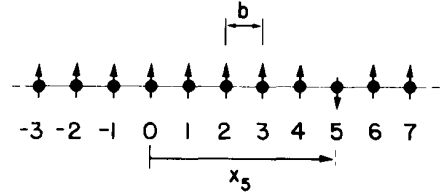


Fig. 15-1. The base state $|x_5\rangle$ of a linear array of spins. All the spins are up except the one at x_5 , which is down.

We can put this trial solution into our usual Hamiltonian equation

$$i\hbar \frac{dC_n}{dt} = \sum_m H_{nm} C_m, \quad (15.10)$$

using Eq (15.8) for the matrix elements. Of course we get an infinite number of equations, but they can all be written as

$$Ea_n = 2Aa_n - Aa_{n-1} - Aa_{n+1} \quad (15.11)$$

We have again exactly the same problem we worked out in Chapter 13, except that where we had E_0 we now have $2A$. The solutions correspond to amplitudes C_n (the down-spin amplitude) which propagate along the lattice with a propagation constant k and an energy

$$E = 2A(1 - \cos kb), \quad (15.12)$$

where b is the lattice constant.

The definite energy solutions correspond to “waves” of down spin—called “spin waves.” And for each wavelength there is a corresponding energy. For large wavelengths (small k) this energy varies as

$$E = Ab^2k^2. \quad (15.13)$$

Just as before, we can consider a localized wave packet (containing, however, only long wavelengths) which corresponds to a spin-down electron in one part of the lattice. This down spin will behave like a “particle.” Because its energy is related to k by (15.13) the “particle” will have an effective mass:

$$m_{\text{eff}} = \frac{\hbar^2}{2Ab^2}. \quad (15.14)$$

These “particles” are sometimes called “magnons.”

15-2 Two spin waves

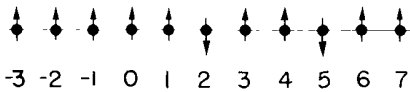


Fig. 15-2. A state with two down spins.

Now we would like to discuss what happens if there are *two* down spins. Again we pick a set of base states. We'll choose states in which there are down spins at two atomic locations, such as the state shown in Fig. 15-2. We can label such a state by the x -coordinates of the two sites with down spins. The one shown can be called $|x_2, x_5\rangle$. In general the base states are $|x_n, x_m\rangle$ —a doubly infinite set! In this system of description, the state $|x_4, x_9\rangle$ and the state $|x_9, x_4\rangle$ are exactly the same state, because each simply says that there is a down spin at 4 and one at 9; there is no meaning to the order. Furthermore, the state $|x_4, x_4\rangle$ has no meaning, there isn't such a thing. We can describe any state $|\psi\rangle$ by giving the amplitudes to be in each of the base states. Thus $C_{m,n} = \langle x_m, x_n | \psi \rangle$ now means the amplitude for a system in the state $|\psi\rangle$ to be in a state in which both the m th and n th atoms have a down spin. The complications which now arise are not complications of ideas—they are merely complexities in bookkeeping. (One of the complexities of quantum mechanics is just the bookkeeping. With more and more down spins, the notation becomes more and more elaborate with lots of indices and the equations always look very horrifying, but the ideas are not necessarily more complicated than in the simplest case.)

The equations of motion of the spin system are the differential equations for the $C_{n,m}$. They are

$$i\hbar \frac{dC_{n,m}}{dt} = \sum_{i,j} (H_{nm,ij}) C_{ij}. \quad (15.15)$$

Suppose we want to find the stationary states. As usual, the derivatives with respect to time become E times the amplitudes and the $C_{m,n}$ can be replaced by the

coefficients $a_{m,n}$. Next we have to work out carefully the effect of H on a state with spins m and n down. It is not hard to figure out. Suppose for a moment that m and n are far enough apart so that we don't have to worry about the obvious trouble. The operation of exchange at the location x_n will move the down spin either to the $(n + 1)$ or $(n - 1)$ atom, and so there's an amplitude that the present state has come from the state $|x_m, x_{n+1}\rangle$ and also an amplitude that it has come from the state $|x_m, x_{n-1}\rangle$. Or it may have been the other spin that moved; so there's a certain amplitude that $C_{m,n}$ is fed from $C_{m+1,n}$ or from $C_{m-1,n}$. These effects should all be equal. The final result for the Hamiltonian equation on $C_{m,n}$ is

$$Ea_{m,n} = -A(a_{m+1,n} + a_{m-1,n} + a_{m,n+1} + a_{m,n-1}) + 4Aa_{m,n}. \quad (15.16)$$

This equation is correct except in two situations. If $m = n$ there is no equation at all, and if $m = n \pm 1$, then two of the terms in Eq. (15.16) should be missing. *We are going to disregard these exceptions.* We simply ignore the fact that some few of these equations are slightly altered. After all, the crystal is supposed to be infinite, and we have an infinite number of terms; neglecting a few might not matter much. So for a first rough approximation let's forget about the altered equations. In other words, we assume that Eq. (15.16) is true for all m and n , even for m and n next to each other. *This is the essential part of our approximation.*

Then the solution is not hard to find. We get immediately

$$C_{m,n} = a_{m,n} e^{-iEt/\hbar}, \quad (15.17)$$

with

$$a_{m,n} = (\text{const.}) e^{ik_1 x_m} e^{ik_2 x_n}, \quad (15.18)$$

where

$$E = 4A - 2A \cos k_1 b - 2A \cos k_2 b. \quad (15.19)$$

Think for a moment what would happen if we had two *independent, single* spin waves (as in the previous section) corresponding to $k = k_1$ and $k = k_2$; they would have energies, from Eq. (15.12), of

$$\epsilon_1 = (2A - 2A \cos k_1 b)$$

and

$$\epsilon_2 = (2A - 2A \cos k_2 b).$$

Notice that the energy E in Eq. (15.19) is just their sum,

$$E = \epsilon(k_1) + \epsilon(k_2). \quad (15.20)$$

In other words we can think of our solution in this way. There are two particles—that is, two spin waves. One of them has a momentum described by k_1 , the other by k_2 , and the energy of the system is the sum of the energies of the two objects. The two particles act completely independently. That's all there is to it.

Of course we have made some approximations, but we do not wish to discuss the precision of our answer at this point. However, you might guess that in a reasonable size crystal with billions of atoms—and, therefore, billions of terms in the Hamiltonian—leaving out a few terms wouldn't make much of an error. If we had so many down spins that there was an appreciable density, then we would certainly have to worry about the corrections.

[Interestingly enough, an exact solution can be written down if there are just the *two* down spins. The result is not particularly important. But it is interesting that the equations can be solved exactly for this case. The solution is:

$$a_{m,n} = \exp^{[ik_c(x_m + x_n)]} \sin k |x_m - x_n|, \quad (15.21)$$

with the energy

$$E = 4A - 2A \cos k_1 b - 2A \cos k_2 b,$$

and with the wave numbers k_c and k related to k_1 and k_2 by

$$k_1 = k_c - k, \quad k_2 = k_c + k. \quad (15.22)$$

This solution includes the “interaction” of the two spins. It describes the fact that when the spins come together there is a certain chance of scattering. The spins act very much like particles with an interaction. But the detailed theory of their scattering goes beyond what we want to talk about here.]

15-3 Independent particles

In the last section we wrote down a Hamiltonian, Eq. (15.15), for a two-particle system. Then using an approximation which is equivalent to neglecting any “interaction” of the two particles, we found the stationary states described by Eqs. (15.17) and (15.18). This state is just the product of two single-particle states. The solution we have given for $a_{m,n}$ in Eq. (15.18) is, however, really not satisfactory. We have very carefully pointed out earlier that the state $|x_9, x_4\rangle$ is *not* a different state from $|x_4, x_9\rangle$ —the *order* of x_m and x_n has no significance. In general, the algebraic expression for the amplitude $C_{m,n}$ must be unchanged if we interchange the values of x_m and x_n , since that doesn’t change the state. Either way, it should represent the amplitude to find a down spin at x_m and a down spin at x_n . But notice that (15.18) is *not* symmetric in x_m and x_n —since k_1 and k_2 can in general be different.

The trouble is that we have not forced our solution of Eq. (15.15) to satisfy this additional condition. Fortunately it is easy to fix things up. Notice first that a solution of the Hamiltonian equation just as good as (15.18) is

$$a_{m,n} = K e^{ik_2 x_m} e^{ik_1 x_n}. \quad (15.23)$$

It even has the same energy we got for (15.18). Any linear combination of (15.15) and (15.23) is also a good solution, and has an energy still given by Eq. (15.19). The solution we should have chosen—because of our symmetry requirement—is just the sum of (15.15) and (15.23):

$$a_{m,n} = K [e^{ik_1 x_m} e^{ik_2 x_n} + e^{ik_2 x_m} e^{ik_1 x_n}]. \quad (15.24)$$

Now, given any k_1 and k_2 the amplitude $C_{m,n}$ is independent of which way we put x_m and x_n —if we should happen to define x_m and x_n reversed we get the same amplitude. Our interpretation of Eq. (15.24) in terms of “magnons” must also be different. We can no longer say that the equation represents *one* particle with wave number k_1 and a *second* particle with wave number k_2 . The amplitude (15.24) represents *one* state with two particles (magnons). The *state* is characterized by the two wave numbers k_1 and k_2 . Our solution looks like a compound state of one particle with the momentum $p_1 = \hbar/k_1$ and another particle with the momentum $p_2 = \hbar/k_2$, but in our state we can’t say which particle is which.

By now, this discussion should remind you of Chapter 4 and our story of identical particles. We have just been showing that the particles of the spin waves—the magnons—behave like identical Bose particles. All amplitudes must be symmetric in the coordinates of the two particles—which is the same as saying that if we “interchange the two particles,” we get back the same amplitude and with the same sign. But, you may be thinking, why did we choose to *add* the two terms in making Eq. (15.24). Why not subtract? With a minus sign, interchanging x_m and x_n would just change the sign of $a_{m,n}$ which doesn’t matter. But interchanging x_m and x_n *doesn’t change anything*—all the electrons of the crystal are exactly where they were before, so there is no reason for even the sign of the amplitude to change. The magnons will behave like Bose particles.†

† In general, the quasi particles of the kind we are discussing may act like either Bose particles or Fermi particles, and as for free particles, the particles with integral spin are bosons and those with half-integral spins are fermions. The “magnon” stands for a spin-up electron turned over. The *change* in spin is *one*. The magnon has an integral spin, and is a boson.

The main points of this discussion have been twofold: First, to show you something about spin waves, and, second, to demonstrate a state whose amplitude is a *product* of two amplitudes, and whose energy is the *sum* of the energies corresponding to the two amplitudes. For *independent particles* the amplitude is the product and the energy is the sum. You can easily see why the energy is the sum. The energy is the coefficient of t in an imaginary exponential—it is proportional to the frequency. If two objects are doing something, one of them with the amplitude $e^{-iE_1 t/\hbar}$ and the other with the amplitude $e^{-iE_2 t/\hbar}$, and if the amplitude for the two things to happen together is the product of the amplitudes for each, then there is a single frequency in the product which is the sum of the two frequencies. The energy corresponding to the amplitude product is the sum of the two energies.

We have gone through a rather long-winded argument to tell you a simple thing. When you don't take into account any interaction between particles, you can think of each particle independently. They can individually exist in the various different states they would have alone, and they will each contribute the energy they would have had if they were alone. However, you must remember that if they are identical particles, they may behave either as Bose or as Fermi particles depending upon the problem. Two extra electrons added to a crystal, for instance, would have to behave like Fermi particles. When the positions of two electrons are interchanged, the amplitude must reverse sign. In the equation corresponding to Eq. (15.24) there would have to be a minus sign between the two terms on the right. As a consequence, two Fermi particles cannot be in exactly the same condition—with equal spins and equal k 's. The amplitude for this state is zero.

15-4 The benzene molecule

Although quantum mechanics provides the basic laws that determine the structures of molecules, these laws can be applied exactly only to the most simple compounds. The chemists have, therefore, worked out various approximate methods for calculating some of the properties of complicated molecules. We would now like to show you how the independent particle approximation is used by the organic chemists. We begin with the benzene molecule.

We discussed the benzene molecule from another point of view in Chapter 10. There we took an approximate picture of the molecule as a two-state system, with the two base states shown in Fig. 15-3. There is a ring of six carbons with a hydrogen bonded to the carbon at each location. With the conventional picture of valence bonds it is necessary to assume double bonds between half of the carbon atoms, and in the lowest energy condition there are the two possibilities shown in the figure. There are also other, higher-energy states. When we discussed benzene in Chapter 10, we just took the two states and forgot all the rest. We found that the ground-state energy of the molecule was not the energy of one of the states in the figure, but was lower than that by an amount proportional to the amplitude to flip from one of these states to the other.

Now we're going to look at the same molecule from a completely different point of view—using a different kind of approximation. The two points of view will give us different answers, but if we improve either approximation it should lead to the truth, a valid description of benzene. However, if we don't bother to improve them, which is of course the usual situation, then you should not be surprised if the two descriptions do not agree exactly. We shall at least show that also with the new point-of-view the lowest energy of the benzene molecule is lower than either of the three-bond structures of Fig. 15-3.

Now we want to use the following picture. Suppose we imagine the six carbon atoms of a benzene molecule connected only by single bonds as in Fig. 15-4. We have removed six electrons—since a bond stands for a pair of electrons—so we have a six-times ionized benzene molecule. Now we will consider what happens when we put back the six electrons one at a time, imagining that each one can run freely around the ring. We assume also that all the bonds shown in Fig. 15-4 are satisfied, and don't need to be considered further.

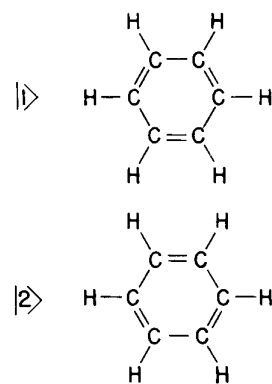


Fig. 15-3. The two base states for the benzene molecule used in Chapter 10.

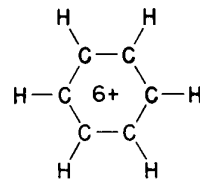


Fig. 15-4. A benzene ring with six electrons removed.

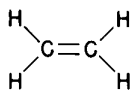


Fig. 15-5. The ethylene molecule.

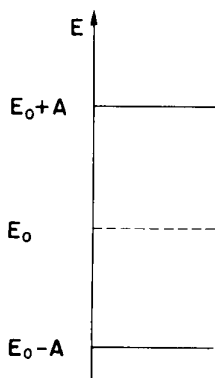


Fig. 15-6. The possible energy levels for the "extra" electrons in the ethylene molecule.

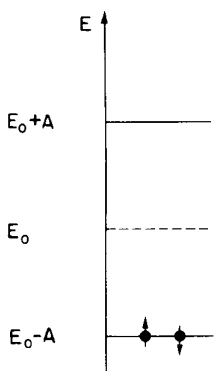


Fig. 15-7. In the extra bond of the ethylene molecule two electrons (one spin up, one spin down) can occupy the lowest energy level.

What happens when we put one electron back into the molecular ion? It might, of course, be located in any one of the six positions around the ring—corresponding to six base states. It would also have a certain amplitude, say A , to go from one position to the next. If we analyze the stationary states, there would be certain possible energy levels. That's only for one electron.

Next put a second electron in. And now we make the most ridiculous approximation that you can think of—that *what one electron does is not affected by what the other is doing*. Of course they really will interact; they repel each other through the Coulomb force, and furthermore when they are both at the same site, they must have considerably different energy than twice the energy for one being there. Certainly the approximation of independent particles is not legitimate when there are only six sites—particularly when we want to put in *six* electrons. Nevertheless the organic chemists have been able to learn a lot by making this kind of an approximation.

Before we work out the benzene molecule in detail, let's consider a simpler example—the ethylene molecule which contains just two carbon atoms with two hydrogen atoms on either side as shown in Fig. 15-5. This molecule has one "extra" bond involving two electrons between the two carbon atoms. Now remove one of these electrons; what do we have? We can look at it as a two-state system—the remaining electron can be at one carbon or the other. We can analyze it as a two-state system. The possible energies for the single electron are either $(E_0 - A)$ or $(E_0 + A)$, as shown in Fig. 15-6.

Now add the second electron. Good, if we have two electrons, we can put the first one in the lower state and the second one in the upper. Not quite; we forgot something. Each one of the states is really double. When we say there's a possible state with the energy $(E_0 - A)$, there are really two. Two electrons can go into the same state if one has its spin up and the other, its spin down. (No more can be put in because of the exclusion principle.) So there really are two possible states of energy $(E_0 - A)$. We can draw a diagram, as in Fig. 15-7, which indicates both the energy levels and their occupancy. In the condition of lowest energy both electrons will be in the lowest state with their spins opposite. The energy of the extra bond in the ethylene molecule therefore is $2(E_0 - A)$ if we neglect the interaction between the two electrons.

Let's go back to the benzene. Each of the two states of Fig. 15-3 has three double bonds. Each of these is just like the bond in ethylene, and contributes $2(E_0 - A)$ to the energy if E_0 is now the energy to put an electron on a site in benzene and A is the amplitude to flip to the next site. So the energy should be roughly $6(E_0 - A)$. But when we studied benzene before, we got that the energy was lower than the energy of the structure with three extra bonds. Let's see if the energy for benzene comes out lower than three bonds from our new point of view.

We start with the six-times ionized benzene ring and add one electron. Now we have a six-state system. We haven't solved such a system yet, but we know what to do. We can write six equations in the six amplitudes, and so on. But let's save some work—by noticing that we've already solved the problem, when we worked out the problem of an electron on an infinite line of atoms. Of course, the benzene is not an infinite line, it has 6 atomic sites in a circle. But imagine that we open out the circle to a line, and number the atoms along the line from 1 to 6. In an infinite line the next location would be 7, but if we insist that this location be identical with number 1 and so on, the situation will be just like the benzene ring. In other words we can take the solution for an infinite line *with an added requirement* that the solution must be periodic with a cycle six atoms long. From Chapter 13 the electron on a line has states of definite energy when the amplitude at each site is $e^{ikx_n} = e^{ikbn}$. For each k the energy is

$$E = E_0 - 2A \cos kb. \quad (15.25)$$

We want to use now only those solutions which repeat every 6 atoms. Let's do first the general case for a ring of N atoms. If the solution is to have a period

of N atomic spacing, e^{ikbN} must be unity; or kbN must be a multiple of 2π . Taking s to represent any integer, our condition is that

$$kbN = 2\pi s. \quad (15.26)$$

We have seen before that there is no meaning to taking k 's outside the range $\pm\pi/b$. This means that we get all possible states by taking values of s in the range $\pm N/2$.

We find then that for an N -atom ring there are N definite energy states† and they have wave numbers k_s given by

$$k_s = \frac{2\pi}{Nb} s. \quad (15.27)$$

Each state has the energy (15.25). We have a line spectrum of possible energy levels. The spectrum for benzene ($N = 6$) is shown in Fig. 15-8(b). (The numbers in parentheses indicate the number of *different* states with the same energy.)

There's a nice way to visualize the six energy levels, as we have shown in part (a) of the figure. Imagine a circle centered on a level with E_0 , and with a radius of $2A$. If we start at the bottom and mark off six equal arcs (at angles from the bottom point of $k_s b = 2\pi s/N$, or $2\pi s/6$ for benzene), then the vertical heights of the points on the circle are the solutions of Eq. (15.25). The six points represent the six possible states. The lowest energy level is at $(E_0 - 2A)$; there are two states with the same energy ($E_0 - A$), and so on.‡ These are possible states for one electron. If we have more than one electron, two—with opposite spins—can go into each state.

For the benzene molecule we have to put in six electrons. For the ground state they will go into the lowest possible energy states—two at $s = 0$, two at $s = +1$, and two at $s = -1$. According to the independent particle approximation the energy of the ground state is

$$\begin{aligned} E_{\text{ground}} &= 2(E_0 - 2A) + 4(E_0 - A) \\ &= 6E_0 - 8A. \end{aligned} \quad (15.28)$$

The energy is indeed less than that of three separate double bonds—by the amount $2A$.

By comparing the energy of benzene to the energy of ethylene it is possible to determine A . It comes out to be 0.8 electron volt, or, in the units the chemists like, 18 kilocalories per mole.

We can use this description to calculate or understand other properties of benzene. For example, using Fig. 15-8 we can discuss the excitation of benzene by light. What would happen if we tried to excite one of the electrons? It could move up to one of the empty higher states. The lowest energy of excitation would be a transition from the highest filled level to the lowest empty level. That takes the energy $2A$. Benzene will absorb light of frequency ν when $h\nu = 2A$. There will also be absorption of photons with the energies $3A$ and $4A$. Needless to say, the absorption spectrum of benzene has been measured and the pattern of spectral lines is more or less correct except that the lowest transition occurs in the ultraviolet; and to fit the data one would have to choose a value of A between 1.4 and 2.4 electron volts. That is, the numerical value of A is two or three times larger than is predicted from the chemical binding energy.

What the chemist does in situations like this is to analyze many molecules of a similar kind and get some empirical rules. He learns, for example: For calculating binding energy use such and such a value of A , but for getting the absorption spectrum approximately right use another value of A . You may feel

† You might think that for N an even number there are $N + 1$ states. That is not so because $s = \pm N/2$ give the same state.

‡ When there are two states (which will have different amplitude distributions) with the same energy, we say that the two states are "degenerate." Notice that *four* electrons can have the energy $E_0 - A$.

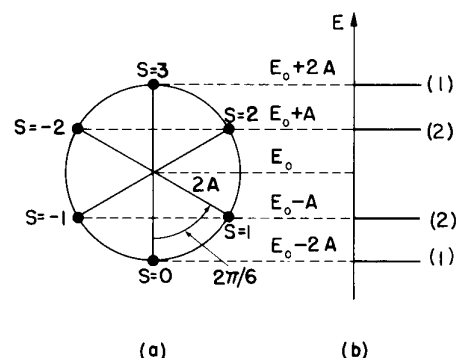


Fig. 15-8. The energy levels in a ring with six electron locations (for example, a benzene ring).

that this sounds a little absurd. It is not very satisfactory from the point of view of a physicist who is trying to understand nature from first principles. But the problem of the chemist is different. He must try to guess ahead of time what is going to happen with molecules that haven't been made yet, or which aren't understood completely. What he needs is a series of empirical rules; it doesn't make much difference where they come from. So he uses the theory in quite a different way than the physicist. He takes equations that have some shadow of the truth in them, but then he must alter the constants in them—making empirical corrections.

In the case of benzene, the principal reason for the inconsistency is our assumption that the electrons are independent—the theory we started with is really not legitimate. Nevertheless, it has some shadow of the truth because its results seem to be going in the right direction. With such equations plus some empirical rules—including various exceptions—the organic chemist makes his way through the morass of complicated things he chooses to study. (Don't forget that the reason a physicist can really calculate from first principles is that he chooses only simple problems. He never solves a problem with 42 or even 6 electrons in it. So far, he has been able to calculate reasonably accurately only the hydrogen atom and the helium atom.)

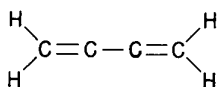


Fig. 15-9. The valence bond representation of the molecule butadiene (1, 3).

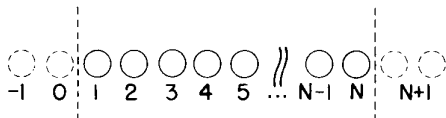


Fig. 15-10. A line of N molecules.

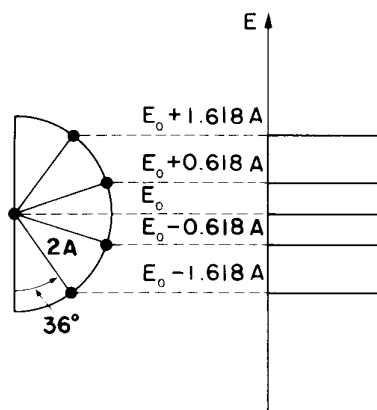


Fig. 15-11. The energy levels of butadiene.

15-5 More organic chemistry

Let's see how the same ideas can be used to study other molecules. Consider a molecule like butadiene (1, 3)—it is drawn in Fig. 15-9 according to the usual valence bond picture.

We can play the same game with the extra four electrons corresponding to the two double bonds. If we remove four electrons, we have four carbon atoms in a line. You already know how to solve a line. You say, "Oh no, I only know how to solve an *infinite* line." But the solutions for the infinite line also include the ones for a finite line. Watch. Let N be the number of atoms on the line and number them from 1 to N as shown in Fig. 15-10. In writing the equations for the amplitude at position 1 you would not have a term feeding from position 0. Similarly, the equation for position N would differ from the one that we used for an infinite line because there would be nothing feeding from position $N + 1$. But suppose that we can obtain a solution for the infinite line which has the following property: the amplitude to be at atom 0 is zero and the amplitude to be at atom $(N + 1)$ is also zero. Then the set of equations for all the locations from 1 to N on the finite line are also satisfied. You might think no such solution exists for the infinite line because our solutions all looked like e^{ikx_n} which has the same absolute value of the amplitude everywhere. But you will remember that the energy depends only on the absolute value of k , so that another solution, which is equally legitimate for the same energy, would be e^{-ikx_n} . And the same is true of any superposition of these two solutions. By subtracting them we can get the solution $\sin kx_n$, which satisfies the requirement that the amplitude be zero at $x = 0$. It still corresponds to the energy $(E_0 - 2A \cos kb)$. Now by a suitable choice for the value of k we can also make the amplitude zero at x_{N+1} . This requires that $(N + 1)kb$ be a multiple of π , or that

$$kb = \frac{\pi}{(N + 1)} s, \quad (15.29)$$

where s is an integer from 1 to N . (We take only positive k 's because each solution contains $+k$ and $-k$; changing the sign of k gives the same state all over again.) For the butadiene molecule, $N = 4$, so there are four states with

$$kb = \pi/5, \quad 2\pi/5, \quad 3\pi/5, \quad \text{and} \quad 4\pi/5. \quad (15.30)$$

We can represent the energy levels using a circle diagram similar to the one for benzene. This time we use a semicircle divided into five equal parts as shown in Fig. 15-11. The point at the bottom corresponds to $s = 0$, which gives no

state at all. The same is true of the point at the top, which corresponds to $s = N + 1$. The remaining 4 points give us four allowed energies. There are four stationary states, which is what we expect having started with four base states. In the circle diagram, the angular intervals are $\pi/5$ or 36 degrees. The lowest energy comes out ($E_0 - 1.618A$). (Ah, what wonders mathematics holds; the golden mean of the Greeks† gives us the lowest energy state of the butadiene molecule according to this theory!)

Now we can calculate the energy of the butadiene molecule when we put in four electrons. With four electrons, we fill up the lowest two levels, each with two electrons of opposite spin. The total energy is

$$E = 2(E_0 - 1.618A) + 2(E_0 - 0.618A) = 4(E_0 - A) - 0.477A. \quad (15.31)$$

This result seems reasonable. The energy is a little lower than for two simple double bonds, but the binding is not so strong as in benzene. Anyway this is the way the chemist analyzes some organic molecules.

The chemist can use not only the energies but the probability amplitudes as well. Knowing the amplitudes for each state, and which states are occupied, he can tell the probability of finding an electron anywhere in the molecule. Those places where the electrons are more likely to be are apt to be reactive in chemical substitutions which require that an electron be shared with some other group of atoms. The other sites are more likely to be reactive in those substitutions which have a tendency to yield an extra electron to the system.

The same ideas we have been using can give us some understanding of a molecule even as complicated as chlorophyll, one version of which is shown in Fig. 15-12. Notice that the double and single bonds we have drawn with heavy lines form a long closed ring with twenty intervals. The extra electrons of the double bonds can run around this ring. Using the independent particle method we can get a whole set of energy levels. There are strong absorption lines from transitions between these levels which lie in the visible part of the spectrum, and give this molecule its strong color. Similar complicated molecules such as the xanthophylls, which make leaves turn red, can be studied in the same way.

There is one more idea which emerges from the application of this kind of theory in organic chemistry. It is probably the most successful or, at least in a certain sense, the most accurate. This idea has to do with the question: In what situations does one get a particularly strong chemical binding? The answer is very interesting. Take the example, first, of benzene, and imagine the sequence of events that occurs as we start with the six-times ionized molecule and add more and more electrons. We would then be thinking of various benzene ions—negative or positive. Suppose we plot the energy of the ion (or neutral molecule) as a function of the number of electrons. If we take $E_0 = 0$ (since we don't know what it is), we get the curve shown in Fig. 15-13. For the first two electrons the slope of the function is a straight line. For each successive group the slope increases, and there is a discontinuity in slope between the groups of electrons. The slope changes when one has just finished filling a set of levels which all have the same energy and must move up to the next higher set of levels for the next electron.

The actual energy of the benzene ion is really quite different from the curve of Fig. 15-13 because of the interactions of the electrons and because of electrostatic energies we have been neglecting. These corrections will, however, vary with n in a rather smooth way. Even if we were to make all these corrections, the resulting energy curve would still have kinks at those values of n which just fill up a particular energy level.

Now consider a very smooth curve that fits the points on the average like the one drawn in Fig. 15-14. We can say that the points *above* this curve have "higher-than-normal" energies, and the points *below* the curve have "lower-than-normal"

† The ratio of the sides of a rectangle which can be divided into a square and a similar rectangle.

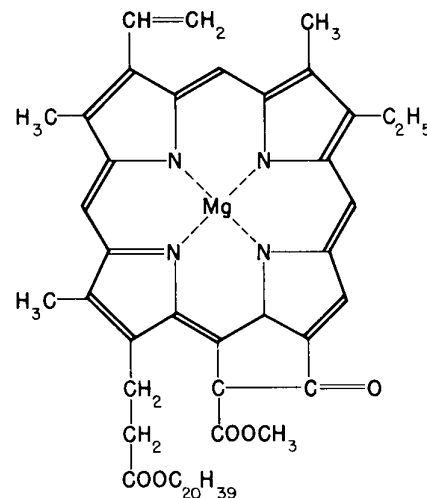


Fig. 15-12. A chlorophyll molecule.

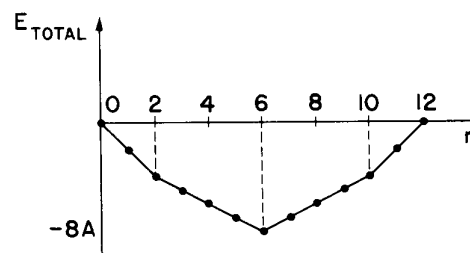


Fig. 15-13. The sum of all the electron energies when the lowest states in Fig. 15-8 are occupied by n electrons if we take that $E_0 = 0$.

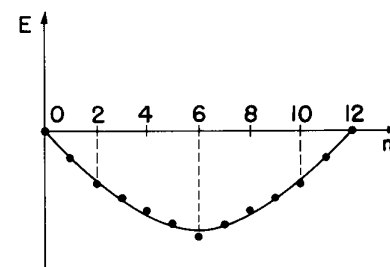


Fig. 15-14. The points of Fig. 15-12 with a smooth curve. Molecules with $n = 2, 6, 10$ are more stable than the others.

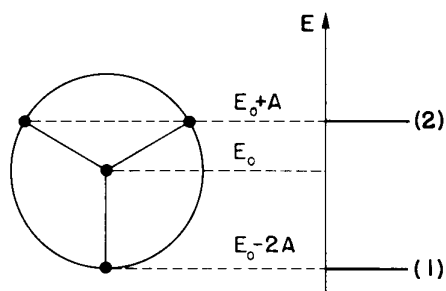


Fig. 15-15. Energy diagram for a ring of three.

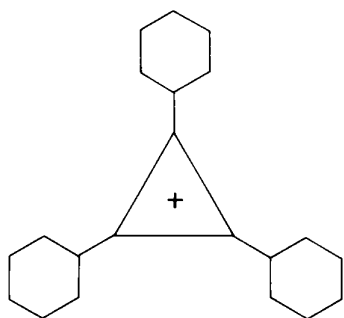


Fig. 15-16. The triphenyl cyclopropenyl cation.

energies. We would, in general, expect that those configurations with a lower-than-normal energy would have an above average stability—chemically speaking. Notice that the configurations farther below the curve always occur at the end of one of the straight line segments—namely when there are enough electrons to fill up an “energy shell,” as it is called. This is the very accurate prediction of the theory. Molecules—or ions—are particularly stable (in comparison with other similar configurations) when the available electrons just fill up an energy shell.

This theory has explained and predicted some very peculiar chemical facts. To take a very simple example, consider a ring of three. It’s almost unbelievable that the chemist can make a ring of three and have it stable, but it has been done. The energy circle for three electrons is shown in Fig. 15-15. Now if you put two electrons in the lower state, you have only two of the three electrons that you require. The third electron must be put in at a much higher level. By our argument this molecule should not be particularly stable, whereas the two-electron structure should be stable. It does turn out, in fact, that the neutral molecule of triphenyl cyclopropenyl is very hard to make, but that the positive ion shown in Fig. 15-16 is relatively easy to make. The ring of three is never really easy because there is always a large stress when the bonds in an organic molecule make an equilateral triangle. To make a stable compound at all, the structure must be stabilized in some way. Anyway if you add three benzene rings on the corners, the positive ion can be made. (The reason for this requirement of added benzene rings is not really understood.)

In a similar way the five-sided ring can also be analyzed. If you draw the energy diagram, you can see in a qualitative way that the six-electron structure should be an especially stable structure, so that such a molecule should be most stable as a negative ion. Now the five-ring is well known and easy to make and always acts as a negative ion. Similarly, you can easily verify that a ring of 4 or 8 is not very interesting, but that a ring of 14 or 10—like a ring of 6—should be especially stable as a neutral object.

15-6 Other uses of the approximation

There are two other similar situations which we will describe only briefly. In considering the structure of an atom, we can consider that the electrons fill successive shells. The Schrödinger theory of electron motion can be worked out easily only for a *single* electron moving in a “central” field—one which varies only with the distance from a point. How can we then understand what goes on in an atom which has 22 electrons?! One way is to use a kind of independent particle approximation. First you calculate what happens with one electron. You get a number of energy levels. You put an electron into the lowest energy state. You can, for a rough model, continue to ignore the electron interactions and go on filling successive shells, but there is a way to get better answers by taking into account—in an approximate way at least—the effect of the electric charge carried by the electron. Each time you add an electron you compute its amplitude to be at various places, and then use this amplitude to estimate a kind of spherically symmetric charge distribution. You use the field of this distribution—together with the field of the positive nucleus and all the previous electrons—to calculate the states available for the next electron. In this way you can get reasonably correct estimates for the energies for the neutral atom and for various ionized states. You find that there are energy shells, just as we saw for the electrons in a ring molecule. With a partially filled shell, the atom will show a preference for taking on one or more extra electrons, or for losing some electrons so as to get into the most stable state of a filled shell.

This theory explains the machinery behind the fundamental chemical properties which show up in the periodic table of the elements. The inert gases are those elements in which a shell has just been completed, and it is especially difficult to make them react. (Some of them do react of course—with fluorine and oxygen, for example; but such compounds are very weakly bound; the so-called inert gases are nearly inert.) An atom which has one electron more or one electron less

than an inert gas will easily lose or gain an electron to get into the especially stable (low-energy) condition which comes from having a completely filled shell—they are the very active chemical elements of valence $+1$ or -1 .

The other situation is found in nuclear physics. In atomic nuclei the protons and neutrons interact with each other quite strongly. Even so, the independent particle model can again be used to analyze nuclear structure. It was first discovered experimentally that nuclei were especially stable if they contained certain particular numbers of neutrons—namely 2, 8, 20, 28, 50, 82. Nuclei containing protons in these numbers are also especially stable. Since there was initially no explanation for these numbers they were called the “magic numbers” of nuclear physics. It is well known that neutrons and protons interact strongly with each other; people were, therefore, quite surprised when it was discovered that an independent particle model predicted a shell structure which came out with the first few magic numbers. The model assumed that each nucleon (proton or neutron) moved in a central potential which was created by the average effects of all the other nucleons. This model failed, however, to give the correct values for the higher magic numbers. Then it was discovered by Maria Mayer, and independently by Jensen and his collaborators, that by taking the independent particle model and adding only a correction for what is called the “spin-orbit interaction,” one could make an improved model which gave all of the magic numbers. (The spin-orbit interaction causes the energy of a nucleon to be lower if its spin has the same direction as its orbital angular momentum from motion in the nucleus.) The theory gives even more—its picture of the so-called “shell structure” of the nuclei enables us to predict certain characteristics of nuclei and of nuclear reactions.

The independent particle approximation has been found useful in a wide range of subjects—from solid-state physics, to chemistry, to biology, to nuclear physics. It is often only a crude approximation, but is able to give an understanding of why there are especially stable conditions—in shells. Since it omits all of the complexity of the interactions between the individual particles, we should not be surprised that it often fails completely to give correctly many important details.

The Dependence of Amplitudes on Position

16-1 Amplitudes on a line

We are now going to discuss how the probability amplitudes of quantum mechanics vary in space. In some of the earlier chapters you may have had a rather uncomfortable feeling that some things were being left out. For example, when we were talking about the ammonia molecule, we chose to describe it in terms of two base states. For one base state we picked the situation in which the nitrogen atom was "above" the plane of the three hydrogen atoms, and for the other base state we picked the condition in which the nitrogen atom was "below" the plane of the three hydrogen atoms. Why did we pick just these two states? Why is it not possible that the nitrogen atom could be at 2 angstroms above the plane of the three hydrogen atoms, or at 3 angstroms, or at 4 angstroms above the plane? Certainly, there are many positions that the nitrogen atom could occupy. Again when we talked about the hydrogen molecular ion, in which there is one electron shared by two protons, we imagined two base states: one for the electron in the neighborhood of proton number one, and the other for the electron in the neighborhood of proton number two. Clearly we were leaving out many details. The electron is not exactly at proton number two but is only in the neighborhood. It could be somewhere above the proton, somewhere below the proton, somewhere to the left of the proton, or somewhere to the right of the proton.

We intentionally avoided discussing these details. We said that we were interested in only certain features of the problem, so we were imagining that when the electron was in the vicinity of proton number one, it would take up a certain rather definite condition. In that condition the probability to find the electron would have some rather definite distribution around the proton, but we were not interested in the details.

We can also put it another way. In our discussion of a hydrogen molecular ion we chose an approximate description when we described the situation in terms of two base states. In reality there are lots and lots of these states. An electron can take up a condition around a proton in its lowest, or ground, state, but there are also many excited states. For each excited state the distribution of the electron around the proton is different. We ignored these excited states, saying that we were interested in only the conditions of low energy. But it is just these other excited states which give the possibility of various distributions of the electron around the proton. If we want to describe in detail the hydrogen molecular ion, we have to take into account also these other possible base states. We could do this in several ways, and one way is to consider in greater detail states in which the location of the electron in space is more carefully described.

We are now ready to consider a more elaborate procedure which will allow us to talk in detail about the position of the electron, by giving a probability amplitude to find the electron anywhere and everywhere in a given situation. This more complete theory provides the underpinning for the approximations we have been making in our earlier discussions. In a sense, our early equations can be derived as a kind of approximation to the more complete theory.

You may be wondering why we did not begin with the more complete theory and make the approximations as we went along. We have felt that it would be much easier for you to gain an understanding of the basic machinery of quantum mechanics by beginning with the two-state approximations and working gradually up to the more complete theory than to approach the subject the other way around. For this reason our approach to the subject appears to be in the reverse order to the one you will find in many books.

16-1 Amplitudes on a line

16-2 The wave function

16-3 States of definite momentum

16-4 Normalization of the states in x

16-5 The Schrödinger equation

16-6 Quantized energy levels

As we go into the subject of this chapter you will notice that we are breaking a rule we have always followed in the past. Whenever we have taken up any subject we have always tried to give a more or less complete description of the physics—showing you as much as we could about where the ideas led to. We have tried to describe the general consequences of a theory as well as describing some specific detail so that you could see where the theory would lead. We are now going to break that rule; we are going to describe how one can talk about probability amplitudes in space and show you the differential equations which they satisfy. We will not have time to go on and discuss many of the obvious implications which come out of the theory. Indeed we will not even be able to get far enough to relate this theory to some of the approximate formulations we have used earlier—for example, to the hydrogen molecule or to the ammonia molecule. For once, we must leave our business unfinished and open-ended. We are approaching the end of our course, and we must satisfy ourselves with trying to give you an introduction to the general ideas and with indicating the connections between what we have been describing and some of the other ways of approaching the subject of quantum mechanics. We hope to give you enough of an idea that you can go off by yourself and by reading books learn about many of the implications of the equations we are going to describe. We must, after all, leave something for the future.

Let's review once more what we have found out about how an electron can move along a line of atoms. When an electron has an amplitude to jump from one atom to the next, there are definite energy states in which the probability amplitude for finding the electron is distributed along the lattice in the form of a traveling wave. For long wavelengths—for small values of the wave number k —the energy of the state is proportional to the square of the wave number. For a crystal lattice with the spacing b , in which the amplitude per unit time for the electron to jump from one atom to the next is iA/\hbar , the energy of the state is related to k (for small kb) by

$$E = Ak^2b^2 \quad (16.1)$$

(see Section 13-3). We also saw that groups of such waves with similar energies would make up a wave packet which would behave like a classical particle with a mass m_{eff} given by:

$$m_{\text{eff}} = \frac{\hbar^2}{2Ab^2} \quad (16.2)$$

Since waves of probability amplitude in a crystal behave like a particle, one might well expect that the general quantum mechanical description of a particle would show the same kind of wave behavior we observed for the lattice. Suppose we were to think of a lattice on a line and imagine that the lattice spacing b were to be made smaller and smaller. In the limit we would be thinking of a case in which the electron could be anywhere along the line. We would have gone over to a continuous distribution of probability amplitudes. We would have the amplitude to find an electron anywhere along the line. This would be one way to describe the motion of an electron in a vacuum. In other words, if we imagine that space can be labeled by an infinity of points all very close together and we can work out the equations that relate the amplitudes at one point to the amplitudes at neighboring points, we will have the quantum mechanical laws of motion of an electron in space.

Let's begin by recalling some of the general principles of quantum mechanics. Suppose we have a particle which can exist in various conditions in a quantum mechanical system. Any particular condition an electron can be found in, we call a "state," which we label with a state vector, say $|\phi\rangle$. Some other condition would be labeled with another state vector, say $|\psi\rangle$. We then introduce the idea of base states. We say that there is a set of states $|1\rangle$, $|2\rangle$, $|3\rangle$, $|4\rangle$, and so on, which have the following properties. First, all of these states are quite distinct—we say they are orthogonal. By this we mean that for any two of the base states $|i\rangle$ and $|j\rangle$ the amplitude $\langle i | j \rangle$ that an electron known to be in the state $|i\rangle$ is also in the

state $|j\rangle$ is equal to zero—unless, of course, $|i\rangle$ and $|j\rangle$ stand for the same state. We represent this symbolically by

$$\langle i | j \rangle = \delta_{ij}. \quad (16.3)$$

You will remember that $\delta_{ij} = 0$ if i and j are different, and $\delta_{ij} = 1$ if i and j are the same number.

Second, the base states $|i\rangle$ must be a complete set, so that any state at all can be described in terms of them. That is, any state $|\phi\rangle$ at all can be described completely by giving all of the amplitudes $\langle i | \phi \rangle$ that a particle in the state $|\phi\rangle$ will also be found in the state $|i\rangle$. In fact, the state vector $|\phi\rangle$ is equal to the sum of the base states each multiplied by a coefficient which is the amplitude of the state $|\phi\rangle$ is also in the state $|i\rangle$:

$$|\phi\rangle = \sum_i |i\rangle \langle i | \phi \rangle. \quad (16.4)$$

Finally, if we consider any two states $|\phi\rangle$ and $|\psi\rangle$, the amplitude that the state $|\psi\rangle$ will also be in the state $|\phi\rangle$ can be found by first projecting the state $|\psi\rangle$ into the base states and then projecting from each base state into the state $|\phi\rangle$. We write that in the following way:

$$\langle \phi | \psi \rangle = \sum_i \langle \phi | i \rangle \langle i | \psi \rangle. \quad (16.5)$$

The summation is, of course, to be carried out over the whole set of base state $|i\rangle$.

In Chapter 13 when we were working out what happens with an electron placed on a linear array of atoms, we chose a set of base states in which the electron was localized at one or other of the atoms in the line. The base state $|n\rangle$ represented the condition in which the electron was localized at atom number “ n .” (There is, of course, no significance to the fact that we called our base states $|n\rangle$ instead of $|i\rangle$.) A little later, we found it convenient to label the base states by the coordinate x_n of the atom rather than by the number of the atom in the array. The state $|x_n\rangle$ is just another way of writing the state $|n\rangle$. Then, following the general rules, any state at all, say $|\psi\rangle$ is described by giving the amplitudes and that an electron in the state $|\psi\rangle$ is also in one of the states $|x_n\rangle$. For convenience we have chosen to let the symbol C_n stand for these amplitudes,

$$C_n = \langle x_n | \psi \rangle. \quad (16.6)$$

Since the base states are associated with a location along the line, we can think of the amplitude C_n as a function of the coordinate x and write it as $C(x_n)$. The amplitudes $C(x_n)$ will, in general, vary with time and are, therefore, also functions of t . We will not generally bother to show explicitly this dependence.

In Chapter 13 we then proposed that the amplitudes $C(x_n)$ should vary with time in a way described by the Hamiltonian equation (Eq. 13.3). In our new notation this equation is

$$i\hbar \frac{\partial C(x_n)}{\partial t} = E_0 C(x_n) - AC(x_n + b) - AC(x_n - b). \quad (16.7)$$

The last two terms on the right-hand side represent the process in which an electron at atom $(n + 1)$ or at atom $(n - 1)$ can feed into atom n .

We found that Eq. (16.7) has solutions corresponding to definite energy states, which we wrote as

$$C(x_n) = e^{iEt/\hbar} e^{ikx_n}. \quad (16.8)$$

For the low-energy states the wavelengths are large (k is small), and the energy is related to k by

$$E = (E_0 - 2A) + Ak^2b^2, \quad (16.9)$$

or, choosing our zero of energy so that $(E_0 - 2A) = 0$, the energy is given by Eq. (16.1).

Let's see what might happen if we were to let the lattice spacing b go to zero, keeping the wave number k fixed. If that is all that were to happen the last term in Eq. (16.9) would just go to zero and there would be no physics. But suppose A and b are varied together so that as b goes to zero the product Ab^2 is kept constant†—using Eq. (16.2) we will write Ab^2 as the constant $\hbar^2/2m_{\text{eff}}$. Under these circumstances, Eq. (16.9) would be unchanged, but what would happen to the differential equation (16.7)?

First we will rewrite Eq. (16.7) as

$$i\hbar \frac{\partial C(x_n)}{\partial t} = (E_0 - 2A)C(x_n) + A[2C(x_n) - C(x_n + b) - C(x_n - b)]. \quad (16.10)$$

For our choice of E_0 , the first term drops out. Next, we can think of a continuous function $C(x)$ that goes smoothly through the proper values $C(x_n)$ at each x_n . As the spacing b goes to zero, the points x_n get closer and closer together, and (if we keep the variation of $C(x)$ fairly smooth) the quantity in the brackets is just proportional to the second derivative of $C(x)$. We can write—as you can see by making a Taylor expansion of each term—the equality

$$2C(x) - C(x + b) - C(x - b) \approx -b^2 \frac{\partial^2 C(x)}{\partial x^2}. \quad (16.11)$$

In the limit, then, as b goes to zero, keeping $b^2 A$ equal to K , Eq. (16.7) goes over into

$$i\hbar \frac{\partial C(x)}{\partial t} = - \frac{\hbar^2}{2m_{\text{eff}}} \frac{\partial^2 C(x)}{\partial x^2}. \quad (16.12)$$

We have an equation which says that the time rate of change of $C(x)$ —the amplitude to find the electron at x —depends on the amplitude to find the electron at nearby points in a way which is proportional to the second derivative of the amplitude with respect to position.

The correct quantum mechanical equation for the motion of an electron in free space was first discovered by Schrödinger. For motion along a line it has exactly the form of Eq. (16.12) if we replace m_{eff} by m , the free-space mass of the electron. For motion along a line in free space the Schrödinger equation is

$$i\hbar \frac{\partial C(x)}{\partial t} = - \frac{\hbar^2}{2m} \frac{\partial^2 C(x)}{\partial x^2}. \quad (16.13)$$

We do not intend to have you think we have derived the Schrödinger equation but only wish to show you one way of thinking about it. When Schrödinger first wrote it down, he gave a kind of derivation based on some heuristic arguments and some brilliant intuitive guesses. Some of the arguments he used were even false, but that does not matter; the only important thing is that the ultimate equation gives a correct description of nature. The purpose of our discussion is then simply to show you that the correct fundamental quantum mechanical equation (16.13) has the same form you get for the limiting case of an electron moving along a line of atoms. This means that we can think of the differential equation in (16.13) as describing the diffusion of a probability amplitude from one point to the next along the line. That is, if an electron has a certain amplitude to be at one point, it will, a little time later, have some amplitude to be at neighboring points. In fact, the equation looks something like the diffusion equations which we have used in Volume I. But there is one main difference: the imaginary coefficient in front of the time derivative makes the behavior completely different from the ordinary diffusion such as you would have for a gas spreading out along a thin tube. Ordinary diffusion gives rise to real exponential solutions, whereas the solutions of Eq. (16.13) are complex waves.

† You can imagine that as the points x_n get closer together, the amplitude A to jump from $x_{n \pm 1}$ to x_n will increase.

16-2 The wave function

Now that you have some idea about how things are going to look, we want to go back to the beginning and study the problem of describing the motion of an electron along a line without having to consider states connected with atoms on a lattice. We want to go back to the beginning and see what ideas we have to use if we want to describe the motion of a free particle in space. Since we are interested in the behavior of a particle along a continuum, we will be dealing with an infinite number of possible states and, as you will see, the ideas we have developed for dealing with a finite number of states will need some technical modifications.

We begin by letting the state vector $|x\rangle$ stand for a state in which a particle is located precisely at the coordinate x . For every value x along the line—for instance 1.73, or 9.67, or 10.00—there is the corresponding state. We will take these states $|x\rangle$ as our base states and, if we include all the points on the line, we will have a complete set for motion in one dimension. Now suppose we have a different kind of a state, say $|\psi\rangle$, in which an electron is distributed in some way along the line. One way of describing this state is to give all the amplitudes that the electron will be also found in each of the base states $|x\rangle$. We must give an infinite set of amplitudes, one for each value of x . We will write these amplitudes as $\langle x | \psi \rangle$. Each of these amplitudes is a complex number and since there is one such complex number for each value of x , the amplitude $\langle x | \psi \rangle$ is indeed just a function of x . We will also write it as $C(x)$,

$$C(x) \equiv \langle x | \psi \rangle. \quad (16.14)$$

We have already considered such amplitudes which vary in a continuous way with the coordinates when we talked about the variations of amplitude with time in Chapter 7. We showed there, for example, that a particle with a definite momentum should be expected to have a particular variation of its amplitude in space. If a particle has a definite momentum p and a corresponding definite energy E , the amplitude to be found at any position x would look like

$$\langle x | \psi \rangle = C(x) \propto e^{+ipx/\hbar}. \quad (16.15)$$

This equation expresses an important general principle of quantum mechanics which connects the base states corresponding to different positions in space to another system of base states—all the states of definite momentum. The definite momentum states are often more convenient than the states in x for certain kinds of problems. Either set of base states is, of course, equally acceptable for a description of a quantum mechanical situation. We will come back later to the matter of the connection between them. For the moment we want to stick to our discussion of a description in terms of the states $|x\rangle$.

Before proceeding, we want to make one small change in notation which we hope will not be too confusing. The function $C(x)$, defined in Eq. (16.14), will of course have a form which depends on the particular state $|\psi\rangle$ under consideration. We should indicate that in some way. We could, for example, specify which function $C(x)$ we are talking about by a subscript say, $C_\psi(x)$. Although this would be a perfectly satisfactory notation, it is a little bit cumbersome and is not the one you will find in most books. Most people simply omit the letter C and use the symbol ψ to define the function

$$\psi(x) \equiv C_\psi(x) = \langle x | \psi \rangle. \quad (16.16)$$

Since this is the notation used by everybody else in the world, you might as well get used to it so that you will not be frightened when you come across it somewhere else. Remember though, that we will now be using ψ in two different ways. In Eq. (16.14), ψ stands for a label we have given to a particular physical state of the electron. On the left-hand side of Eq. (16.16), on the other hand, the symbol ψ is used to define a mathematical function of x which is equal to the amplitude to be associated with each point x along the line. We hope it will not be too confusing

once you get accustomed to the idea. Incidentally, the function $\psi(x)$ is usually called “the wave function”—because it more often than not has the form of a complex wave in its variables.

Since we have defined $\psi(x)$ to be the amplitude that an electron in the state ψ will be found at the location x , we would like to interpret the absolute square of ψ to be the probability of finding an electron at the position x . Unfortunately, the probability of finding a particle exactly at any particular point is zero. The electron will, in general, be smeared out in a certain region of the line, and since, in any small piece of the line, there are an infinite number of points, the probability that it will be at any one of them cannot be a finite number. We can only describe the probability of finding an electron in terms of a *probability distribution*[†] which gives the *relative* probability of finding the electron at various approximate locations along the line. Let's let $\text{prob}(x, \Delta x)$ stand for the chance of finding the electron in a small interval Δx located near x . If we go to a small enough scale in any physical situation, the probability will be varying smoothly from place to place, and the probability of finding the electron in any small finite line segment Δx will be proportional to Δx . We can modify our definitions to take this into account.

We can think of the amplitude $\langle x | \psi \rangle$ as representing a kind of “amplitude density” for all the base states $|x\rangle$ in a small region. Since the probability of finding an electron in a small interval Δx at x should be proportional to the interval Δx , we choose our definition of $\langle x | \psi \rangle$ so that the following relation holds:

$$\text{prob}(x, \Delta x) = |\langle x | \psi \rangle|^2 \Delta x.$$

The amplitude $\langle x | \psi \rangle$ is therefore proportional to the amplitude that an electron in the state ψ will be found in the base state x and the constant of proportionality is chosen so that the absolute square of the amplitude $\langle x | \psi \rangle$ gives the *probability density* of finding an electron in any small region. We can write, equivalently,

$$\text{prob}(x, \Delta x) = |\psi(x)|^2 \Delta x. \quad (16.17)$$

We will now have to modify some of our earlier equations to make them compatible with this new definition of a probability amplitude. Suppose we have an electron in the state $|\psi\rangle$ and we want to know the amplitude for finding it in a different state $|\phi\rangle$ which may correspond to a different spread-out condition of the electron. When we were talking about a finite set of discrete states, we would have used Eq. (16.5). Before modifying our definition of the amplitudes we would have written

$$\langle \phi | \psi \rangle = \sum_{\text{all } x} \langle \phi | x \rangle \langle x | \psi \rangle. \quad (16.18)$$

Now if both of these amplitudes are normalized in the same way as we have described above, then a sum of all the states in a small region of x would be equivalent to multiplying by Δx , and the sum over all values of x simply becomes an integral. With our modified definitions, the correct form becomes

$$\langle \phi | \psi \rangle = \int_{\text{all } x} \langle \phi | x \rangle \langle x | \psi \rangle dx. \quad (16.19)$$

The amplitude $\langle x | \psi \rangle$ is what we are now calling $\psi(x)$ and, in a similar way, we will choose to let the amplitude $\langle x | \phi \rangle$ be represented by $\phi(x)$. Remembering that $\langle \phi | x \rangle$ is the complex conjugate of $\langle x | \phi \rangle$, we can write Eq. (16.18) as

$$\langle \phi | \psi \rangle = \int \phi^*(x) \psi(x) dx. \quad (16.20)$$

With our new definitions everything follows with the same formulas as before if you always replace a summation sign by an integral over x .

We should mention one qualification to what we have been saying. Any suitable set of base states must be complete if it is to be used for an adequate

[†] For a discussion of probability distributions see Vol. I, Section 6-4.

description of what is going on. For an electron in one dimension it is not really sufficient to specify only the base states $|x\rangle$, because for each of these states the electron may have a spin which is either up or down. One way of getting a complete set is to take two sets of states in x , one for up spin and the other for down spin. We will, however, not worry about such complications for the time being.

16-3 States of definite momentum

Suppose we have an electron in a state $|\psi\rangle$ which is described by the probability amplitude $\langle x|\psi\rangle = \psi(x)$. We know that this represents a state in which the electron is spread out along the line in a certain distribution so that the probability of finding the electron in a small interval dx at the location x is just

$$\text{prob}(x, dx) = |\psi(x)|^2 dx.$$

What can we say about the momentum of this electron? We might ask what is the probability that this electron has the momentum p ? Let's start out by calculating the amplitude that the state $|\psi\rangle$ is in another state $|\text{mom } p\rangle$ which we define to be a state with the definite momentum p . We can find this amplitude by using our basic equation for the resolution of amplitudes, Eq. (16.20). In terms of the state $|\text{mom } p\rangle$

$$\langle \text{mom } p | \psi \rangle = \int_{x=-\infty}^{+\infty} \langle \text{mom } p | x \rangle \langle x | \psi \rangle dx. \quad (16.21)$$

And the probability that the electron will be found with the momentum p should be given in terms of the absolute square of this amplitude. We have again, however, a small problem about the normalizations. In general we can only ask about the probability of finding an electron with a momentum in a small range dp at the momentum p . The probability that the momentum is exactly some value p must be zero (unless the state $|\psi\rangle$ happens to be a state of definite momentum). Only if we ask for the probability of finding the momentum in a small range dp at the momentum p will we get a finite probability. There are several ways the normalizations can be adjusted. We will choose one of them which we think to be the most convenient, although that may not be apparent to you just now.

We take our normalizations so that the probability is related to the amplitude by

$$\text{prob}(p, dp) = |\langle \text{mom } p | \psi \rangle|^2 \frac{dp}{2\pi\hbar}. \quad (16.22)$$

With this definition the normalization of the amplitude $\langle \text{mom } p | x \rangle$ is determined. The amplitude $\langle \text{mom } p | x \rangle$ is, of course, just the complex conjugate of the amplitude $\langle x | \text{mom } p \rangle$, which is just the one we have written down in Eq. (16.15). With the normalization we have chosen, it turns out that the proper constant of proportionality in front of the exponential is just 1. Namely,

$$\langle \text{mom } p | x \rangle = \langle x | \text{mom } p \rangle^* = e^{-ipx/\hbar}. \quad (16.23)$$

Equation (16.21) then becomes

$$\langle \text{mom } p | \psi \rangle = \int_{-\infty}^{+\infty} e^{-ipx/\hbar} \langle x | \psi \rangle dx. \quad (16.24)$$

This equation together with Eq. (16.22) allows us to find the momentum distribution for any state $|\psi\rangle$.

Let's look at a particular example—for instance one in which an electron is localized in a certain region around $x = 0$. Suppose we take a wave function which has the following form:

$$\psi(x) = K e^{-x^2/4\sigma^2}. \quad (16.25)$$

The probability distribution in x for this wave function is the absolute square, or

$$\text{prob}(x, dx) = P(x) dx = K^2 e^{-x^2/2\sigma^2} dx. \quad (16.26)$$

The probability density function $P(x)$ is the Gaussian curve shown in Fig. 16-1. Most of the probability is concentrated between $x = +\sigma$ and $x = -\sigma$. We say that the “half-width” of the curve is σ . (More precisely, σ is equal to the root-mean-square of the coordinate x for something spread out according to this distribution.) We would normally choose the constant K so that the probability density $P(x)$ is not merely *proportional* to the probability per unit length in x of finding the electron, but has a scale such that $P(x) \Delta x$ is *equal* to the probability of finding the electron in Δx near x . The constant K which does this can be found by requiring that $\int_{-\infty}^{+\infty} P(x) dx = 1$, since there must be unit probability that the electron is found somewhere. Here, we get that $K = (2\pi\sigma^2)^{-1/4}$. [We have used the fact that $\int_{-\infty}^{+\infty} e^{-t^2} dt = \sqrt{\pi}$; see Vol. I, page 40-6.]

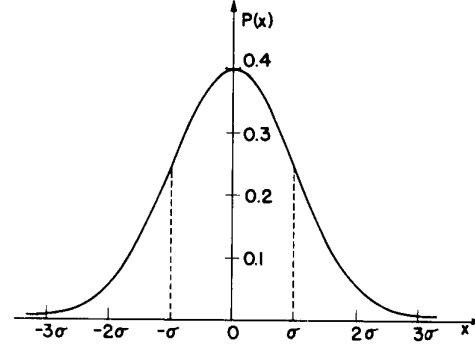


Fig. 16-1. The probability density for the wave function of Eq. (16.24).

Now let's find the distribution in momentum. Let's let $\phi(p)$ stand for the amplitude to find the electron with the momentum p ,

$$\phi(p) \equiv \langle \text{mom } p | \psi \rangle. \quad (16.27)$$

Substituting Eq. (16.25) into Eq. (16.24) we get

$$\phi(p) = \int_{-\infty}^{+\infty} e^{-ipx/\hbar} \cdot K e^{-x^2/4\sigma^2} dx. \quad (16.28)$$

the integral can also be rewritten as

$$K e^{-p^2\sigma^2/\hbar^2} \int_{-\infty}^{+\infty} e^{-(1/4\sigma^2)(x + 2ip\sigma^2/\hbar)^2} dx. \quad (16.29)$$

We can now make the substitution $u = x + 2ip\sigma^2/\hbar$, and the integral is

$$\int_{-\infty}^{+\infty} e^{-u^2/4\sigma^2} du = 2\sigma\sqrt{\pi}. \quad (16.30)$$

(The mathematicians would probably object to the way we got there, but the result is, nevertheless, correct.)

$$\phi(p) = (8\pi\sigma^2)^{1/4} e^{-p^2\sigma^2/\hbar^2}. \quad (16.31)$$

We have the interesting result that the amplitude function in p has precisely the same mathematical form as the amplitude function in x ; only the width of the Gaussian is different. We can write this as

$$\phi(p) = (2\pi\eta^2)^{-1/4} e^{-p^2/4\eta^2}, \quad (16.32)$$

where the half-width η of the p -distribution function is related to the half-width σ of the x -distribution by

$$\eta = \frac{\hbar}{2\sigma}. \quad (16.33)$$

Our result says: if we make the width of the distribution in x very small by making σ small, η becomes large and the distribution in p is very much spread out. Or, conversely: if we have a narrow distribution in p , it must correspond to a spread-out distribution in x . We can, if we like, consider η and σ to be some measure of the uncertainty in the localization of the momentum and of the position of the electron in the state we are studying. If we call them Δp and Δx respectively Eq. (16.33) becomes

$$\Delta p \Delta x = \frac{\hbar}{2}. \quad (16.34)$$

Interestingly enough, it is possible to prove that for any other form of a distribution in x or in p , the product $\Delta p \Delta x$ cannot be smaller than the one we have found here. The Gaussian distribution gives the smallest possible value for the product of the root-mean-square widths. In general, we can say

$$\Delta p \Delta x \geq \frac{\hbar}{2}. \quad (16.35)$$

This is a quantitative statement of the Heisenberg uncertainty principle, which we have discussed qualitatively many times before. We have usually made the approximate statement that the minimum value of the product $\Delta p \Delta x$ is of the same order as \hbar .

16-4 Normalization of the states in x

We return now to the discussion of the modifications of our basic equations which are required when we are dealing with a continuum of base states. When we have a finite number of discrete states, a fundamental condition which must be satisfied by the set of base states is

$$\langle i | j \rangle = \delta_{ij}. \quad (16.36)$$

If a particle is in one base state, the amplitude to be in another base state is 0. By choosing a suitable normalization, we have defined the amplitude $\langle i | i \rangle$ to be 1. These two conditions are described by Eq. (16.36). We want now to see how this relation must be modified when we use the base states $|x\rangle$ of a particle on a line. If the particle is known to be in one of the base states $|x\rangle$, what is the amplitude that it will be in another base state $|x'\rangle$? If x and x' are two different locations along the line, then the amplitude $\langle x | x' \rangle$ is certainly 0, so that is consistent with Eq. (16.36). But if x and x' are equal, the amplitude $\langle x | x' \rangle$ will not be 1, because of the same old normalization problem. To see how we have to patch things up, we go back to Eq. (16.19), and apply this equation to the special case in which the state $|\phi\rangle$ is just the base state $|x'\rangle$. We would have then

$$\langle x' | \psi \rangle = \int \langle x' | x \rangle \psi(x) dx. \quad (16.37)$$

Now the amplitude $\langle x | \psi \rangle$ is just what we have been calling the function $\psi(x)$. Similarly the amplitude $\langle x' | \psi \rangle$, since it refers to the same state ψ , is the same function of the variable x' , namely $\psi(x')$. We can, therefore, rewrite Eq. (16.37) as

$$\psi(x') = \int \langle x' | x \rangle \psi(x) dx. \quad (16.38)$$

This equation must be true for any state ψ and, therefore, for any arbitrary function $\psi(x)$. This requirement should completely determine the nature of the amplitude $\langle x | x' \rangle$ —which is, of course, just a function that depends on x and x' .

Our problem now is to find a function $f(x, x')$ which when multiplied into $\psi(x)$, and integrated over all x gives just the quantity $\psi(x')$. It turns out that there is no mathematical function which will do this! At least nothing like what we ordinarily mean by a “function.”

Suppose we pick x' to be the special number 0 and define the amplitude $\langle 0 | x \rangle$ to be some function of x , let's say $f(x)$. Then Eq. (16.36) would read as follows:

$$\psi(0) = \int f(x)\psi(x) dx. \quad (16.39)$$

What kind of function $f(x)$ could possibly satisfy this equation? Since the integral must not depend on what values $\psi(x)$ takes for values of x other than 0, $f(x)$ must clearly be 0 for all values of x except 0. But if $f(x)$ is 0 everywhere, the integral will be 0, too, and Eq. (16.39) will not be satisfied. So we have an impossible situation: we wish a function to be 0 everywhere but at a point, and still to give a finite integral. Since we can't find a function that does this, the easiest way out is just to say that the function $f(x)$ is *defined* by Eq. (16.37). Namely, $f(x)$ is that function which makes (16.39) correct. The function which does this was first invented by Dirac and carries his name. We write it $\delta(x)$. All we are saying is that the function $\delta(x)$ has the strange property that if it is substituted for $f(x)$ in the Eq. (16.39), the integral picks out the value that $\psi(x)$ takes on when x is equal 0; and, since the integral must be independent of $\psi(x)$ for all values of x other than 0, the function $\delta(x)$ must be 0 everywhere except at $x = 0$. Summarizing, we write

$$\langle 0 | x \rangle = \delta(x), \quad (16.40)$$

where $\delta(x)$ is defined by

$$\psi(0) = \int \delta(x)\psi(x) dx. \quad (16.41)$$

Notice what happens if we use the special function "1" for the function ψ in Eq. (16.41). Then we have the result

$$1 = \int \delta(x) dx. \quad (16.42)$$

That is, the function $\delta(x)$ has the property that it is 0 everywhere except at $x = 0$ but has a finite integral equal to unity. We must imagine that the function $\delta(x)$ has such a fantastic infinity at one point that the total area comes out equal to one.

One way of imagining what the Dirac δ -function is like is to think of a sequence of rectangles—or any other peaked function you care to—which gets narrower and narrower and higher and higher, always keeping a unit area, as sketched in Fig. 16-2. The integral of this function from $-\infty$ to $+\infty$ is always 1. If you multiply it by any function $\psi(x)$ and integrate the product, you get something which is approximately the value of the function at $x = 0$, the approximation getting better and better as you use the narrower and narrower rectangles. You can if you wish, imagine the δ -function in terms of this kind of limiting process. The only important thing, however, is that the δ -function is defined so that Eq. (16.41) is true for every possible function $\psi(x)$. That uniquely defines the δ -function. Its properties are then as we have described.

If we change the argument of the δ -function from x to $x - x'$, the corresponding relations are

$$\begin{aligned} \delta(x - x') &= 0, & x' \neq x, \\ \int \delta(x - x')\psi(x) dx &= \psi(x'). \end{aligned} \quad (16.43)$$

If we use $\delta(x - x')$ for the amplitude $\langle x | x' \rangle$ in Eq. (16.38), that equation is satisfied. Our result then is that for our base states in x , the condition corresponding to (16.36) is

$$\langle x' | x \rangle = \delta(x - x'). \quad (16.44)$$

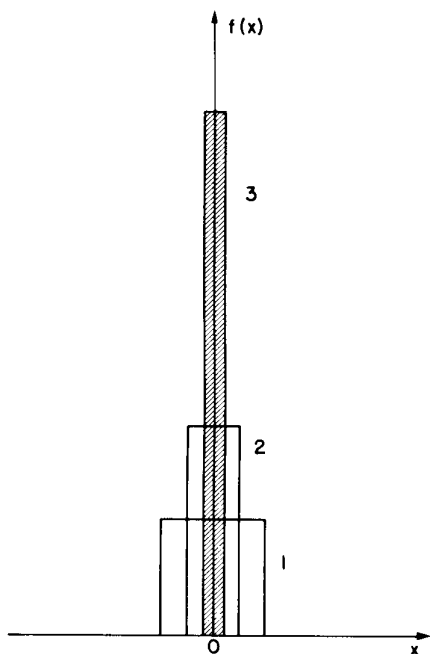


Fig. 16-2. A set of functions, all of unit area, which look more and more like $\delta(x)$.

We have now completed the necessary modifications of our basic equations which are necessary for dealing with the continuum of base states corresponding to the points along a line. The extension to three dimensions is fairly obvious; first we replace the coordinate x by the vector \mathbf{r} . Then integrals over x become replaced by integrals over x, y , and z . In other words, they become volume integrals. Finally, the one-dimensional δ -function must be replaced by just the product of three δ -functions, one in x , one in y , and the other in z , $\delta(x - x') \delta(y - y') \delta(z - z')$. Putting everything together we get the following set of equations for the amplitudes for particle in three dimensions:

$$\langle \phi | \psi \rangle = \int \langle \phi | \mathbf{r} \rangle \langle \mathbf{r} | \psi \rangle d\text{Vol}, \quad (16.45)$$

$$\langle \mathbf{r} | \psi \rangle = \psi(\mathbf{r}), \quad (16.46)$$

$$\langle \mathbf{r} | \phi \rangle = \phi(\mathbf{r}),$$

$$\langle \phi | \psi \rangle = \int \phi^*(\mathbf{r}) \psi(\mathbf{r}) d\text{Vol}, \quad (16.47)$$

$$\langle \mathbf{r}' | \mathbf{r} \rangle = \delta(x - x') \delta(y - y') \delta(z - z'), \quad (16.48)$$

What happens when there is more than one particle? We will tell you about how to handle two particles and you will easily see what you must do if you want to deal with a larger number. Suppose there are two particles, which we can call particle No. 1 and particle No. 2. What shall we use for the base states? One perfectly good set can be described by saying that particle 1 is at x_1 and particle 2 is at x_2 , which we can write as $|x_1 x_2\rangle$. Notice that describing the position of *only one particle does not define* a base state. Each base state must define the condition of the entire system. You must not think that each particle moves independently as a wave in three dimensions. Any physical state $|\psi\rangle$ can be defined by giving all of the amplitudes $\langle x_1, x_2 | \psi \rangle$ to find the two particles at x_1 and x_2 . This generalized amplitude is therefore a function of the *two* sets of coordinates x_1 and x_2 . You see that such a function is not a wave in the sense of an oscillation that moves along in three dimensions. Neither is it generally simply a product of two individual waves, one for each particle. It is, in general, some kind of a wave in the six dimensions defined by x_1 and x_2 . If there are two particles in nature which are interacting, there is no way of describing what happens to one of the particles by trying to write down a wave function for it alone. The famous paradoxes that we considered in earlier chapters—where the measurements made on one particle were claimed to be able to tell what was going to happen to another particle, or were able to destroy an interference—have caused people all sorts of trouble because they have tried to think of the wave function of one particle alone, rather than the correct wave function in the coordinates of both particles. The complete description can be given correctly only in terms of functions of the coordinates of both particles.

16-5 The Schrödinger equation

So far we have just been worrying about how we can describe states which may involve an electron being anywhere at all in space. Now we have to worry about putting into our description the physics of what can happen in various circumstances. As before, we have to worry about how states can change with time. If we have a state $|\psi\rangle$ which goes over into another state $|\psi'\rangle$ sometime later, we can describe the situation for all times by making the wave function—which is just the amplitude $\langle \mathbf{r} | \psi \rangle$ —a function of time as well as a function of the coordinate. A particle in a given situation can then be described by giving a time-varying wave function $\psi(\mathbf{r}, t) = \psi(x, y, z, t)$. This time-varying wave function describes the evolution of successive states that occur as time develops. This so-called “coordinate representation”—which gives the projections of the state $|\psi\rangle$ into the base states $|\mathbf{r}\rangle$ may not always be the most convenient one to use—but we will consider it first.

In Chapter 8 we described how states varied in time in terms of the Hamiltonian H_{ij} . We saw that the time variation of the various amplitudes was given in terms of the matrix equation

$$i\hbar \frac{dC_i}{dt} = \sum_j H_{ij}C_j. \quad (16.49)$$

This equation says that the time variation of each amplitude C_i is proportional to all of the other amplitudes C_j , with the coefficients H_{ij} .

How would we expect Eq. (16.49) to look when we are using the continuum of base states $|x\rangle$? Let's first remember that Eq. (16.49) can also be written as

$$i\hbar \frac{d}{dt} \langle i | \psi \rangle = \sum_j \langle i | \hat{H} | j \rangle \langle j | \psi \rangle.$$

Now it is clear what we should do. For the x -representation we would expect

$$i\hbar \frac{d}{dt} \langle x | \psi \rangle = \int \langle x | \hat{H} | x' \rangle \langle x' | \psi \rangle dx'. \quad (16.50)$$

The sum over the base states $|j\rangle$, gets replaced by an integral over x' . Since $\langle x | \hat{H} | x' \rangle$ should be some function of x and x' , we can write it as $H(x, x')$ —which corresponds to H_{ij} in Eq. (16.49). Then Eq. (16.50) is the same as

$$i\hbar \frac{d}{dt} \psi(x) = \int H(x, x') \psi(x') dx' \quad (16.51)$$

with

$$H(x, x') \equiv \langle x | \hat{H} | x' \rangle.$$

According to Eq. (16.51), the rate of change of the ψ at x would depend on the value of ψ at all other points x' ; the factor $H(x, x')$ is the amplitude per unit time that the electron will jump from x' to x . *It turns out in nature, however, that this amplitude is zero except for points x' very close to x .* This means—as we saw in the example of the chain of atoms at the beginning of the chapter, Eq. (16.12)—that the right-hand side of Eq. (16.15) can be expressed completely in terms of ψ and the derivatives of ψ with respect to x , all evaluated at the position x .

For a particle moving freely in space with no forces, no disturbances, the correct law of physics is

$$\int H(x, x') \psi(x') dx' = -\frac{\hbar^2}{2m} \frac{d^2}{dx^2} \psi(x).$$

Where did we get that from? Nowhere. It's not possible to derive it from anything you know. It came out of the mind of Schrödinger, invented in his struggle to find an understanding of the experimental observations of the real world. You can perhaps get some clue of why it should be that way by thinking of our derivation of Eq. (16.12) which came from looking at the propagation of an electron in a crystal.

Of course, free particles are not very exciting. What happens if we put forces on the particle? Well, if the force of a particle can be described in terms of a scalar potential $V(x)$ —which means we are thinking of electric forces but not magnetic forces—and if we stick to low energies so that we can ignore complexities which come from relativistic motions, then the Hamiltonian which fits the real world gives

$$\int H(x, x') \psi(x') dx' = -\frac{\hbar^2}{2m} \frac{d^2}{dx^2} \psi(x) + V(x) \psi(x). \quad (16.52)$$

Again, you can get some clue as to the origin of this equation if you go back to the motion of an electron in a crystal, and see how the equations would have to be modified if the energy of the electron varied slowly from one atomic site to the other—as it might do if there were an electric field across the crystal. Then

the term E_0 in Eq. (16.7) would vary slowly with position and would correspond to the new term we have added in (16.52).

[You may be wondering why we went straight from Eq. (16.51) to Eq. (16.52) instead of just giving you the correct function for the amplitude $H(x, x') = \langle x | \hat{H} | x' \rangle$. We did that because $H(x, x')$ can only be written in terms of strange algebraic functions, although the whole integral on the right-hand side of Eq. (16.51) comes out in terms of things you are used to. If you are really curious, $H(x, x')$ can be written in the following way:

$$H(x, x') = -\frac{\hbar^2}{2m} \delta''(x - x') + V(x) \delta(x - x'),$$

where δ'' means the second derivative of the delta function. This rather strange function can be replaced by a somewhat more convenient algebraic differential operator, which is completely equivalent:

$$H(x, x') = \left\{ -\frac{\hbar^2}{2m} \frac{d^2}{dx^2} + V(x) \right\} \delta(x - x').$$

We will *not* be using these forms, but will work directly with the form in Eq. (16.52).]

If we now use the expression we have in (16.52) for the integral in (16.50) we get the following differential equation for $\psi(x) = \langle x | \psi \rangle$:

$$i\hbar \frac{\partial \psi}{\partial t} = -\frac{\hbar^2}{2m} \frac{\partial^2}{\partial x^2} \psi(x) + V(x) \psi(x). \quad (16.53)$$

It is fairly obvious what we should use instead of Eq. (16.53) if we are interested in motion in three dimensions. The only changes are that d^2/dx^2 gets replaced by

$$\nabla^2 = \frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} + \frac{\partial^2}{\partial z^2},$$

and $V(x)$ gets replaced by $V(x, y, z)$. The amplitude $\psi(x, y, z)$ for an electron moving in a potential $V(x, y, z)$ obeys the differential equation

$$i\hbar \frac{\partial \psi}{\partial t} = -\frac{\hbar^2}{2m} \nabla^2 \psi + V \psi. \quad (16.54)$$

It is called the Schrödinger equation, and was the first quantum-mechanical equation ever known. It was written down by Schrödinger before any of the other quantum equations we have described in this book were discovered.

Although we have approached the subject along a completely different route, the great historical moment marking the birth of the quantum mechanical description of matter occurred when Schrödinger first wrote down his equation in 1926. For many years the internal atomic structure of matter had been a great mystery. No one had been able to understand what held matter together, why there was chemical binding, and especially how it could be that atoms could be stable. Although Bohr had been able to give a description of the internal motion of an electron in a hydrogen atom which seemed to explain the observed spectrum of light emitted by this atom, the reason that electrons moved in this way remained a mystery. Schrödinger's discovery of the proper equations of motion for electrons on an atomic scale provided a theory from which atomic phenomena could be calculated quantitatively, accurately, and in detail. In principle, Schrödinger's equation is capable of explaining all atomic phenomena except those involving magnetism and relativity. It explains the energy levels of an atom, and all the facts of chemical binding. This is, however, true only in principle—the mathematics soon becomes too complicated to solve exactly any but the simplest problems. Only the hydrogen and helium atoms have been calculated to a high accuracy. However, with various approximations, some fairly sloppy, many of the facts of more complicated atoms and of the chemical binding of molecules can be understood. We have shown you some of these approximations in earlier chapters.

The Schrödinger equation as we have written it does not take into account any magnetic effects. It is possible to take such effects into account in an approximate way by adding some more terms to the equation. However, as we have seen in Volume II, magnetism is essentially a relativistic effect, and so a correct description of the motion of an electron in an arbitrary electromagnetic field can only be discussed in a proper relativistic equation. The correct relativistic equation for the motion of an electron was discovered by Dirac a year after Schrödinger brought forth his equation, and takes on quite a different form. We will not be able to discuss it at all here.

Before we go on to look at some of the consequences of the Schrödinger equation, we would like to show you what it looks like for a system with a large number of particles. We will not be making any use of the equation, but just want to show it to you to emphasize that the wave function ψ is not simply an ordinary wave in space, but is a function of many variables. If there are many particles, the equation becomes

$$-i\hbar \frac{\partial \psi(r_1, r_2, r_3, \dots)}{\partial t} = \sum_i \frac{\hbar^2}{2m_i} \left\{ \frac{\partial^2 \psi}{\partial x_i^2} + \frac{\partial^2 \psi}{\partial y_i^2} + \frac{\partial^2 \psi}{\partial z_i^2} \right\} + V(r_1, r_2, \dots) \psi. \quad (16.55)$$

The potential function V is what corresponds classically to the total potential energy of all the particles. If there are no external forces acting on the particles, the function V is simply the electrostatic energy of interaction of all the particles. That is, if the i th particle carries the charge $Z_i q_e$, then the function V is simply†

$$V(r_1, r_2, r_3, \dots) = \sum_{\substack{\text{all} \\ \text{pairs}}} \frac{Z_i Z_j}{r_{ij}} e^2. \quad (16.56)$$

16-6 Quantized energy levels

In a later chapter we will look in detail at a solution of Schrödinger's equation for a particular example. We would like now, however, to show you how one of the most remarkable consequence of Schrödinger's equation comes about—namely, the surprising fact that a differential equation involving only continuous functions of continuous variables in space can give rise to quantum effects such as the discrete energy levels in an atom. The essential fact to understand is how it can be that an electron which is confined to a certain region of space by some kind of a potential “well” must necessarily have only one or another of a certain well-defined set of discrete energies.

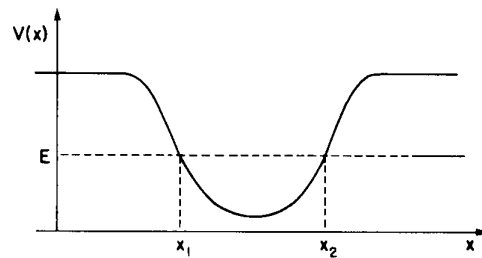


Fig. 16-3. A potential well for a particle moving along x .

Suppose we think of an electron in a one-dimensional situation in which its potential energy varies with x in a way described by the graph in Fig. 16-3. We will assume that this potential is static—it doesn't vary with time. As we have done so many times before, we would like to look for solutions corresponding to states of definite energy, which means, of definite frequency. Let's try a solution of the form

$$\psi = a(x)e^{-iEt/\hbar}. \quad (16.57)$$

† We are using the convention of the earlier volumes according to which $e^2 \equiv q_e^2/4\pi\epsilon_0$.

If we substitute this function into the Schrödinger equation, we find that the function $a(x)$ must satisfy the following differential equation:

$$\frac{d^2 a(x)}{dx^2} = \frac{2m}{\hbar} [V(x) - E]a(x). \quad (16.58)$$

This equation says that at each x the second derivative of $a(x)$ with respect to x is proportional to $a(x)$, the coefficient of proportionality being given by the quantity $(V - E)$. The second derivative of $a(x)$ is the rate of change of its slope. If the potential V is greater than the energy E of the particle, the rate of change of the slope of $a(x)$ will have the same sign as $a(x)$. That means that the curve of $a(x)$ will be concave away from the axis. That is, it will have, more or less, the character of the positive or negative exponential function, $e^{\pm x}$. This means that in the region to the left of x_1 , in Fig. 16-3, where V is greater than the assumed energy E , the function $a(x)$ would have to look like one or another of the curves

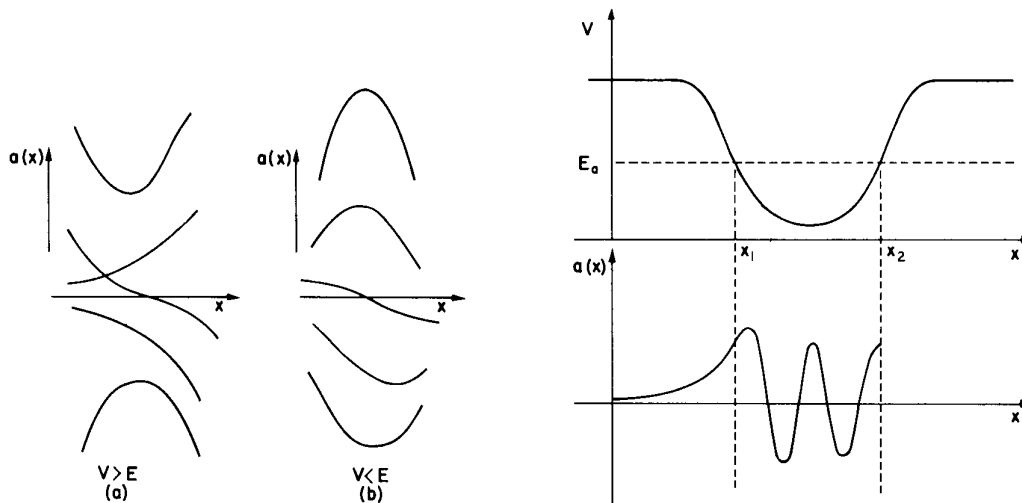


Fig. 16-4. Possible shapes of the wave function $a(x)$ for $V > E$ and for $V < E$.

Fig. 16-5. A wave function for the energy E_a which goes to zero for negative x .

shown in part (a) of Fig. 16-4.

If, on the other hand, the potential function V is less than the energy E , the second derivative of $a(x)$ with respect to x has the opposite sign from $a(x)$ itself, and the curve of $a(x)$ will always be concave toward the axis like one of the pieces shown in part (b) of Fig. 16-4. The solution in such a region has, piece-by-piece, roughly the form of a sinusoidal curve.

Now let's see if we can construct graphically a solution for the function $a(x)$ which corresponds to a particle of energy E_a in the potential V shown in Fig. 16-3. Since we are trying to describe a situation in which a particle is bound *inside* the potential well, we want to look for solutions in which the wave amplitude takes on very small values when x is way outside the potential well. We can easily imagine a curve like the one shown in Fig. 16-5 which tends toward zero for large negative values of x , and grows smoothly as it approaches x_1 . Since V is equal to E_a at x_1 , the curvature of the function becomes zero at this point. Between x_1 and x_2 , the quantity $V - E_a$ is always a negative number, so the function $a(x)$ is always concave toward the axis, and the curvature is larger the larger the difference between E_a and V . If we continue the curve into the region between x_1 and x_2 , it should go more or less as shown in Fig. 16-5.

Now let's continue this curve into the region to the right of x_2 . There it curves away from the axis and takes off toward large positive values, as drawn in Fig. 16-6. For the energy E_a we have chosen, the solution for $a(x)$ gets larger and larger with increasing x . In fact, its *curvature* is also increasing (if the potential continues to stay flat). The amplitude rapidly grows to immense proportions.

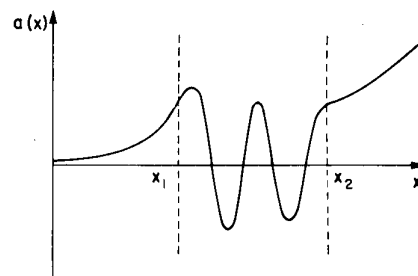


Fig. 16-6. The wave function $a(x)$ of Fig. 16-5 continued beyond x_2 .

Fig. 16-7. The wave function $a(x)$ for an energy E_b greater than E_a .

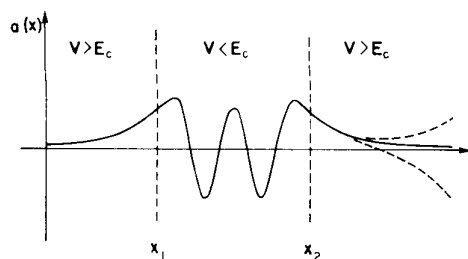
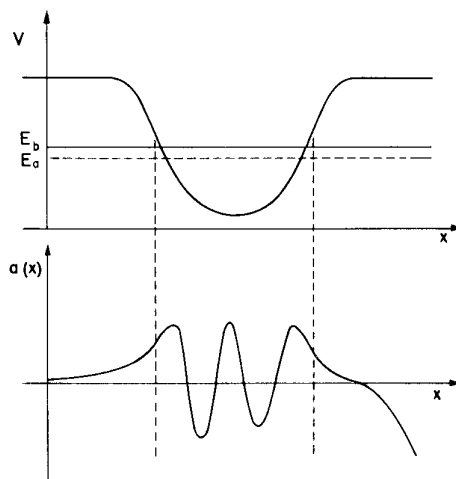
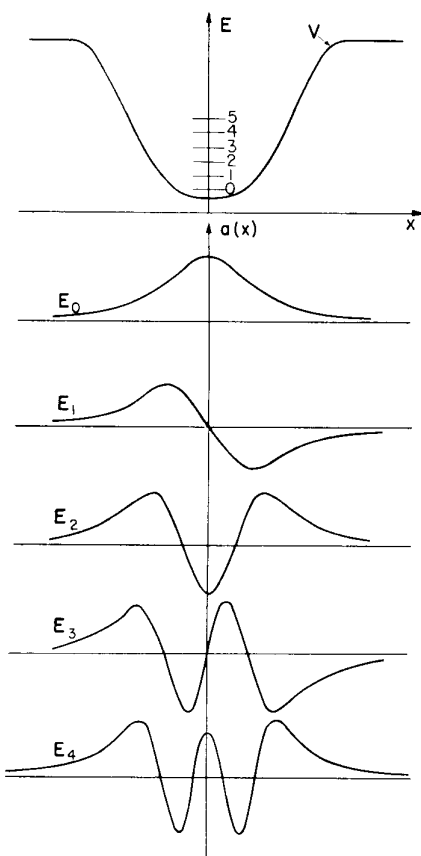


Fig. 16-8. A wave function for the energy E_c between E_a and E_b .



What does this mean? It simply means that the particle is not “bound” in the potential well. It is infinitely more likely to be found *outside* of the well, than *inside*. For the solution we have manufactured, the electron is more likely to be found at $x = +\infty$ than anywhere else. We have failed to find a solution for a bound particle.

Let’s try another energy, say one a little bit higher than E_a —say the energy E_b in Fig. 16-7. If we start with the same conditions on the left, we get the solution drawn in the lower half of Fig. 16-7. It looked at first as though it were going to be better, but it ends up just as bad as the solution for E_a —except that now $a(x)$ is getting more and more *negative* as we go toward large values of x .

Maybe that’s the clue. Since changing the energy a little bit from E_a to E_b causes the curve to flip from one side of the axis to the other, perhaps there is some energy lying between E_a and E_b for which the curve will approach zero for large values of x . There is, indeed, and we have sketched how the solution might look in Fig. 16-8.

You should appreciate that the solution we have drawn in the figure is a very special one. If we were to raise or lower the energy ever so slightly, the function would go over into curves like one or the other of the two broken-line curves shown in Fig. 16-8, and we would not have the proper conditions for a bound particle. We have obtained a result that if a particle is to be bound in a potential well, it can do so only if it has a very definite energy.

Does that mean that there is only one energy for a particle bound in a potential well? No. Other energies are possible, but not energies too close to E_c . Notice that the wave function we have drawn in Fig. 16-8 crosses the axis four times in the region between x_1 and x_2 . If we were to pick an energy quite a bit lower than E_c , we could have a solution which crosses the axis only three times, only two times, only once, or not at all. The possible solutions are sketched in Fig. 16-9. (There may also be other solutions corresponding to values of the energy higher than the ones shown.) Our conclusion is that if a particle is bound in a potential well, its energy can take on only the certain special values in a discrete energy spectrum. You see how a differential equation can describe the basic fact of quantum physics.

We might remark one other thing. If the energy E is above the top of the potential well, then there are no longer any discrete solutions, and any possible energy is permitted. Such solutions correspond to the scattering of free particles by a potential well. We have seen an example of such solutions when we considered the effects of impurity atoms in a crystal.

Fig. 16-9. The function $a(x)$ for the five lowest energy bound states.

Symmetry and Conservation Laws

17-1 Symmetry

In classical physics there are a number of quantities which are *conserved*—such a momentum, energy, and angular momentum. Conservation theorems about corresponding quantities also exist in quantum mechanics. The most beautiful thing of quantum mechanics is that the conservation theorems can, in a sense, be derived from something else, whereas in classical mechanics they are practically the starting points of the laws. (There are ways in classical mechanics to do an analogous thing to what we will do in quantum mechanics, but it can be done only at a very advanced level.) In quantum mechanics, however, the conservation laws are very deeply related to the principle of superposition of amplitudes, and to the symmetry of physical systems under various changes. This is the subject of the present chapter. Although we will apply these ideas mostly to the conservation of angular momentum, the essential point is that the theorems about the conservation of all kinds of quantities are—in the quantum mechanics—related to the symmetries of the system.

We begin, therefore, by studying the question of symmetries of systems. A very simple example is the hydrogen molecular ion—we could equally well take the ammonia molecule—in which there are two states. For the hydrogen molecular ion we took as our base states one in which the electron was located near proton number 1, and another in which the electron was located near proton number 2. The two states—which we called $|1\rangle$ and $|2\rangle$ —are shown again in Fig. 17-1(a). Now, so long as the two nuclei are both exactly the same, then there is a certain *symmetry* in this physical system. That is to say, if we were to *reflect* the system in the plane halfway between the two protons—by which we mean that everything on one side of the plane gets moved to the symmetric position on the other side—we would get the situations in Fig. 17-1(b). Since the protons are identical, the *operation of reflection* changes $|1\rangle$ into $|2\rangle$ and $|2\rangle$ into $|1\rangle$. We'll call this reflection operation \hat{P} and write

$$\hat{P}|1\rangle = |2\rangle, \quad \hat{P}|2\rangle = |1\rangle. \quad (17.1)$$

So our \hat{P} is an operator in the sense that it “*does something*” to a state to make a new state. The interesting thing is that \hat{P} operating on *any* state produces some *other* state of the system.

Now \hat{P} , like any of the other operators we have described, has matrix elements which can be defined by the usual obvious notation. Namely,

$$P_{11} = \langle 1 | \hat{P} | 1 \rangle \quad \text{and} \quad P_{12} = \langle 1 | \hat{P} | 2 \rangle$$

are the matrix elements we get if we multiply $\hat{P}|1\rangle$ and $\hat{P}|2\rangle$ on the left by $\langle 1|$. From Eq. (17.1) they are

$$\begin{aligned} \langle 1 | \hat{P} | 1 \rangle &= P_{11} = \langle 1 | 2 \rangle = 0, \\ \langle 1 | \hat{P} | 2 \rangle &= P_{12} = \langle 1 | 1 \rangle = 1. \end{aligned} \quad (17.2)$$

In the same way we can get P_{21} and P_{22} . The matrix of \hat{P} —with respect to the base system $|1\rangle$ and $|2\rangle$ —is

$$P = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}. \quad (17.3)$$

We see once again that the words *operator* and *matrix* in quantum mechanics are

17-1 Symmetry

17-2 Symmetry and conservation

17-3 The conservation laws

17-4 Polarized light

17-5 The disintegration of the Λ^0

17-6 Summary of the rotation matrices

Review: Chapter 52, Vol. I, *Symmetry in Physical Laws*

Reference: *Angular Momentum in Quantum Mechanics*:
A. R. Edmonds, Princeton University Press, 1957

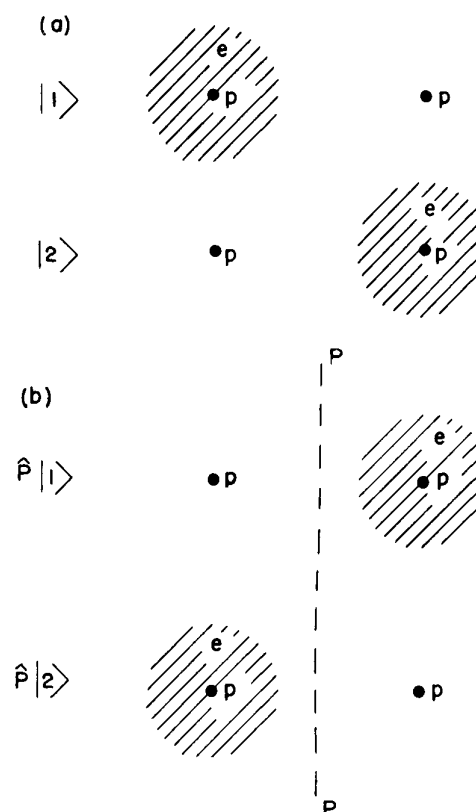


Fig. 17-1. If the states $|1\rangle$ and $|2\rangle$ are reflected in the plane P - P , they go into $|2\rangle$ and $|1\rangle$, respectively.

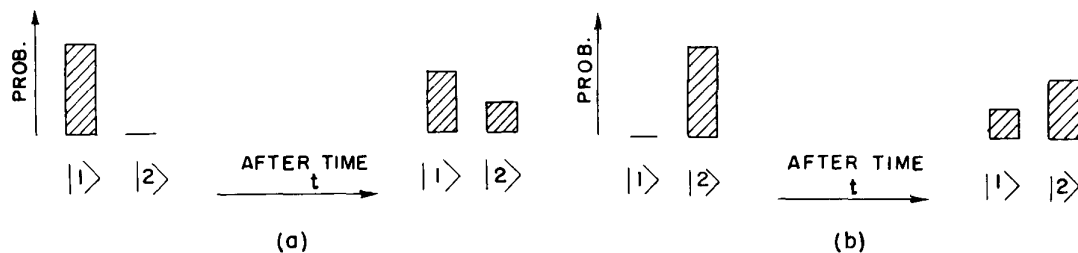


Fig. 17-2. In a symmetric system, if a pure $|1\rangle$ state develops as shown in part (a), a pure $|2\rangle$ state will develop as in part (b).

practically interchangeable. There are slight technical differences—like the difference between a “numeral” and a “number”—but the distinction is something pedantic that we don’t have to worry about. So whether \hat{P} defines an operation, or is actually used to define a matrix of numbers, we will call it interchangeably an operator or a matrix.

Now we would like to point out something. We will *suppose* that the *physics* of the whole hydrogen molecular ion system is *symmetrical*. It doesn’t have to be—it depends, for instance, on what else is near it. But if the system is symmetrical, the following idea should certainly be true. Suppose we start at $t = 0$ with the system in the state $|1\rangle$ and find after an interval of time t that the system turns out to be in a more complicated situation—in some linear combination of the two base states. Remember that in Chapter 8 we used to represent “going for a period of time” by multiplying by the operator \hat{U} . That means that the system would after a while—say 15 seconds to be definite—be in some other state. For example, it might be $\sqrt{2/3}$ parts of the state $|1\rangle$ and $i\sqrt{1/3}$ parts of the state $|2\rangle$, and we would write

$$|\psi \text{ at 15 sec}\rangle = \hat{U}(15, 0) |1\rangle = \sqrt{2/3} |1\rangle + i\sqrt{1/3} |2\rangle. \quad (17.4)$$

Now we ask what happens if we start the system in the *symmetric* state $|2\rangle$ and wait for 15 seconds under the *same conditions*? It is clear that if the world is symmetric—as we are supposing—we should get the state symmetric to (17.4):

$$|\psi \text{ at 15 sec}\rangle = \hat{U}(15, 0) |2\rangle = \sqrt{2/3} |2\rangle + i\sqrt{1/3} |1\rangle. \quad (17.5)$$

The same ideas are sketched diagrammatically in Fig. 17-2. So if the *physics* of a system is symmetrical with respect to some plane, and we work out the behavior of a particular state, we also know the behavior of the state we would get by reflecting the original state in the symmetry plane.

We would like to say the same things a little bit more generally—which means a little more abstractly. Let \hat{Q} be any one of a number of operations that you could perform on a system *without changing the physics*. For instance, for \hat{Q} we might be thinking of \hat{P} , the operation of a *reflection* in the plane between the two atoms in the hydrogen molecule. Or, in a system with two electrons, we might be thinking of the operation of *interchanging* the two electrons. Another possibility would be, in a spherically symmetric system, the operation of a *rotation* of the whole system through a finite angle around some axis—which wouldn’t change the physics. Of course, we would normally want to give each special case some special notation for \hat{Q} . Specifically, we will normally define the $\hat{R}_y(\theta)$ to be the operation “rotate the system about the y -axis by the angle θ ”. By \hat{Q} we mean just any one of the operators we have described or any other one—which leaves the basic physical situation unchanged.

Let’s think of some more examples. If we have an atom with *no external magnetic field* or *no external electric field*, and if we were to turn the coordinates around any axis, it would be the same physical system. Again, the ammonia molecule is symmetrical with respect to a reflection in a plane parallel to that of the three hydrogens—so long as there is no electric field. When there is an electric field, when we make a reflection we would have to change the electric field also,

and that changes the physical problem. But if we have no external field, the molecule is symmetrical.

Now we consider a general situation. Suppose we start with the state $|\psi_1\rangle$ and after some time or other under given physical conditions it has become the state $|\psi_2\rangle$. We can write

$$|\psi_2\rangle = \hat{U} |\psi_1\rangle. \quad (17.6)$$

[You can be thinking of Eq. (17.4).] Now imagine we perform the operation \hat{Q} on the whole system. The state $|\psi_1\rangle$ will be transformed to a state $|\psi'_1\rangle$, which we can also write as $\hat{Q} |\psi_1\rangle$. Also the state $|\psi_2\rangle$ is changed into $|\psi'_2\rangle = \hat{Q} |\psi_2\rangle$. Now if the physics is symmetrical under \hat{Q} (don't forget the *if*; it is not a general property of systems), then, waiting for the same time under the same conditions, we should have

$$|\psi'_2\rangle = \hat{U} |\psi'_1\rangle. \quad (17.7)$$

[Like Eq. (17.5).] But we can write $\hat{Q} |\psi_1\rangle$ for $|\psi'_1\rangle$ and $\hat{Q} |\psi_2\rangle$ for $|\psi'_2\rangle$ so (17.7) can also be written

$$\hat{Q} |\psi_2\rangle = \hat{U} \hat{Q} |\psi_1\rangle. \quad (17.8)$$

If we now replace $|\psi_2\rangle$ by $\hat{U} |\psi_1\rangle$ —Eq. (17.6)—we get

$$\hat{Q} \hat{U} |\psi_1\rangle = \hat{U} \hat{Q} |\psi_1\rangle. \quad (17.9)$$

It's not hard to understand what this means. Thinking of the hydrogen ion it says that: “making a reflection and waiting a while”—the expression on the right of Eq. (17.9)—is the same as “waiting a while and then making a reflection”—the expression on the left of (17.9). These should be the same so long as U doesn't change under the reflection.

Since (17.9) is true for *any* starting state $|\psi_1\rangle$, it is really an equation about the operators:

$$\hat{Q} \hat{U} = \hat{U} \hat{Q}. \quad (17.10)$$

This is what we wanted to get—it is a *mathematical statement of symmetry*. When Eq. (17.10) is true, we say that the operators \hat{U} and \hat{Q} *commute*. We can then *define* “symmetry” in the following way: A physical system is *symmetric* with respect to the operation \hat{Q} when \hat{Q} commutes with \hat{U} , the operation of the passage of time. [In terms of matrices, the product of two operators is equivalent to the matrix product, so Eq. (17.10) also holds for the matrices Q and U for a system which is symmetric under the transformation Q .]

Incidentally, since for infinitesimal times ϵ we have $\hat{U} = 1 - i\hat{H}\epsilon/\hbar$ —where \hat{H} is the usual Hamiltonian (see Chapter 8)—you can see that if (17.10) is true, it is also true that

$$\hat{Q} \hat{H} = \hat{H} \hat{Q}. \quad (17.11)$$

So (17.11) is the mathematical statement of the condition for the symmetry of a physical situation under the operator \hat{Q} . It *defines* a symmetry.

17-2 Symmetry and conservation

Before applying the result we have just found, we would like to discuss the idea of symmetry a little more. Suppose that we have a very special situation: after we operate on a state with \hat{Q} , we get the same state. This is a very special case, but let's suppose it happens to be true for a state $|\psi_0\rangle$ that $|\psi'\rangle = \hat{Q} |\psi_0\rangle$ is physically the same state as $|\psi_0\rangle$. That means that $|\psi'\rangle$ is equal to $|\psi_0\rangle$ except for some phase factor.[†] How can that happen? For instance, suppose that we

[†] Incidentally, you can show that \hat{Q} is necessarily a *unitary operator*—which means that if it operates on $|\psi\rangle$ to give some number times $|\psi\rangle$, the number must be of the form $e^{i\delta}$, where δ is real. It's a small point, and the proof rests on the following observation. Any operation like a reflection or a rotation doesn't lose any particles, so the normalization of $|\psi'\rangle$ and $|\psi\rangle$ must be the same; they can only differ by a pure imaginary phase factor.

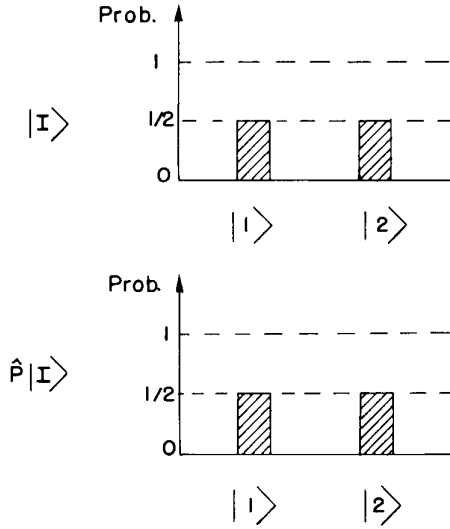


Fig. 17-3. The state $|I\rangle$ and the state $\hat{P}|I\rangle$ obtained by reflecting $|I\rangle$ in the central plane.

have an H_2^+ ion in the state which we once called $|I\rangle$. For this state there is equal amplitude to be in the base states $|1\rangle$ and $|2\rangle$. The probabilities are shown as a bar graph in Fig. 17-3(a). If we operate on $|I\rangle$ with the reflection operator \hat{P} , it flips the state over changing $|1\rangle$ to $|2\rangle$ and $|2\rangle$ to $|1\rangle$ —we get the probabilities shown in Fig. 17-3(b). But that's just the state $|I\rangle$ all over again. If we start with state $|II\rangle$ the probabilities before and after reflection look just the same. However, there is a difference if we look at the *amplitudes*. For the state $|I\rangle$ the amplitudes are the *same* after the reflection, but for the state $|II\rangle$ the amplitudes have the opposite sign. In other words,

$$\begin{aligned}\hat{P}|I\rangle &= \hat{P}\left\{\frac{|1\rangle + |2\rangle}{\sqrt{2}}\right\} = \frac{|2\rangle + |1\rangle}{\sqrt{2}} = |I\rangle, \\ \hat{P}|II\rangle &= \hat{P}\left\{\frac{|1\rangle - |2\rangle}{\sqrt{2}}\right\} = \frac{|2\rangle - |1\rangle}{\sqrt{2}} = -|II\rangle.\end{aligned}\quad (17.12)$$

If we write $\hat{P}|\psi_0\rangle = e^{i\delta}|\psi_0\rangle$, we have that $e^{i\delta} = 1$ for the state $|I\rangle$ and $e^{i\delta} = -1$ for the state $|II\rangle$.

Let's look at another example. Suppose we have a RHC polarized photon propagating in the z -direction. If we do the operation of a rotation around the z -axis, we know that this just multiplies the amplitude by $e^{i\phi}$ when ϕ is the angle of the rotation. So for the rotation operation in this case, δ is just equal to the angle of rotation.

Now it is clear that *if it happens to be true* that an operator \hat{Q} just changes the phase of a state at some time, say $t = 0$, *it is true forever*. In other words, if the state $|\psi_1\rangle$ goes over into the state $|\psi_2\rangle$ after a time t , or

$$\hat{U}(t, 0)|\psi_1\rangle = |\psi_2\rangle \quad (17.13)$$

and if the symmetry of the situation makes it so that

$$\hat{Q}|\psi_1\rangle = e^{i\delta}|\psi_1\rangle, \quad (17.14)$$

then it is also true that

$$\hat{Q}|\psi_2\rangle = e^{i\delta}|\psi_2\rangle. \quad (17.15)$$

This is clear, since

$$\hat{Q}|\psi_2\rangle = \hat{Q}\hat{U}|\psi_1\rangle = \hat{U}\hat{Q}|\psi_1\rangle,$$

and if $\hat{Q}|\psi_1\rangle = e^{i\delta}|\psi_1\rangle$, then

$$\hat{Q}|\psi_2\rangle = \hat{U}e^{i\delta}|\psi_1\rangle = e^{i\delta}\hat{U}|\psi_1\rangle = e^{i\delta}|\psi_2\rangle.$$

[The sequence of equalities follows from (17.13) and (17.10) for a symmetrical system, from (17.14), and from the fact that a number like $e^{i\delta}$ commutes with an operator.]

So with certain symmetries something which is true initially is true for all times. But isn't that just a *conservation law*? Yes! It says that if you look at the original state and by making a little computation on the side discover that an operation *which is a symmetry operation of the system* produces only a multiplication by a certain phase, then you know that the same property will be true of the final state—the same operation multiplies the final state by the same phase factor. This is always true even though we may not know anything else about the inner mechanism of the universe which changes a system from the initial to the final state. Even if we do not care to look at the details of the machinery by which the system gets from one state to another, we can still say that if a thing is in a state with a certain symmetry character originally, and if the Hamiltonian for this thing is symmetrical under that symmetry operation, then the state will have the same symmetry character for all times. That's the basis of all the conservation laws of quantum mechanics.

Let's look at a special example. Let's go back to the \hat{P} operator. We would like first to modify a little our definition of \hat{P} . We want to take for \hat{P} not just a

mirror reflection, because that requires defining the plane in which we put the mirror. There is a special kind of a reflection that doesn't require the specification of a plane. Suppose we redefine the operation \hat{P} this way: First you reflect in a mirror in the z -plane so that z goes to $-z$, x stays x , and y stays y ; then you turn the system 180° about the z -axis so that x is made to go to $-x$ and y to $-y$. The whole thing is called an *inversion*. Every point is projected *through the origin* to the diametrically opposite position. All the coordinates of everything are reversed. We will still use the symbol \hat{P} for this operation. It is shown in Fig. 17-4. It is a little more convenient than a simple reflection because it doesn't require that you specify which coordinate plane you used for the reflection—you need specify only the point which is at the center of symmetry.

Now let's suppose that we have a state $|\psi_0\rangle$ which under the inversion operation goes into $e^{i\delta}|\psi_0\rangle$ —that is,

$$|\psi'_0\rangle = \hat{P}|\psi_0\rangle = e^{i\delta}|\psi_0\rangle. \quad (17.16)$$

Then suppose that we invert again. After *two* inversions we are right back where we started from—nothing is changed at all. We must have that

$$\hat{P}|\psi'_0\rangle = \hat{P}\hat{P}|\psi_0\rangle = |\psi_0\rangle.$$

But

$$\hat{P}\hat{P}|\psi_0\rangle = \hat{P}e^{i\delta}|\psi_0\rangle = e^{i\delta}\hat{P}|\psi_0\rangle = (e^{i\delta})^2|\psi_0\rangle.$$

It follows that

$$(e^{i\delta})^2 = 1.$$

So *if the inversion operator is a symmetry operation* of a state, there are only two possibilities for δ :

$$e^{i\delta} = \pm 1,$$

which means that

$$\hat{P}|\psi_0\rangle = |\psi_0\rangle \quad \text{or} \quad \hat{P}|\psi_0\rangle = -|\psi_0\rangle. \quad (17.17)$$

Classically, if a state is symmetric under an inversion, the operation gives back the same state. In quantum mechanics, however, there are the two possibilities: we get the *same state* or *minus* the same state. When we get the *same state*, $\hat{P}|\psi_0\rangle = |\psi_0\rangle$, we say that the state $|\psi_0\rangle$ has *even parity*. When the sign is reversed so that $\hat{P}|\psi_0\rangle = -|\psi_0\rangle$, we say that the state has *odd parity*. (The inversion operator \hat{P} is also known as the parity operator.) The state $|I\rangle$ of the H_2^+ ion has even parity; and the state $|II\rangle$ has odd parity—see Eq. (17.12). There are, of course, states which are not symmetric under the operation \hat{P} ; these are states with no definite parity. For instance, in the H_2^+ system the state $|I\rangle$ has even parity, the state $|II\rangle$ has odd parity, and the state $|I\rangle$ has no definite parity.

When we speak of an operation like inversion being performed “*on a physical system*” we can think about it in two ways. We can think of *physically moving* whatever is at \mathbf{r} to the inverse point at $-\mathbf{r}$, or we can think of *looking* at the same system from a new frame of reference x', y', z' related to the old by $x' = -x$, $y' = -y$, and $z' = -z$. Similarly, when we think of rotations, we can think of rotating bodily a physical system, or of rotating the coordinate frame with respect to which we measure the system, keeping the “system” fixed in space. Generally, the two points of view are essentially equivalent. For rotation they are equivalent *except* that rotating a *system* by the angle θ is like rotating the reference frame by the *negative* of θ . In these lectures we have usually considered what happens when a projection is made into a new set of axes. What you get that way is the same as what you get if you leave the axes fixed and rotate the system *backwards* by the same amount. When you do that, the signs of the angles are reversed.†

† In other books you may find formulas with different signs; they are probably using a different definition of the angles.

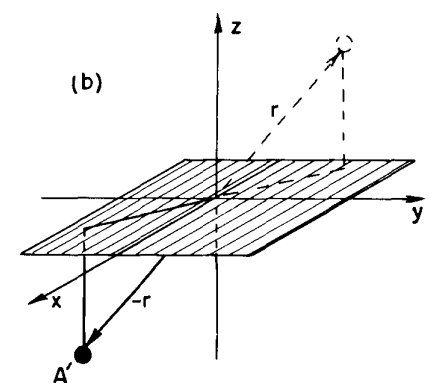
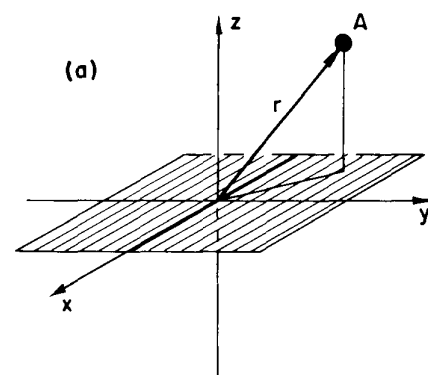


Fig. 17-4. The operation of inversion, \hat{P} . Whatever is at the point A at (x, y, z) is moved to the point A' at $(-x, -y, -z)$.

Many of the *laws* of physics—but not all—are unchanged by a reflection or an inversion of the coordinates. They are *symmetric* with respect to an inversion. The laws of electrodynamics, for instance, are unchanged if we change x to $-x$, y to $-y$, and z to $-z$ in *all* the equations. The same is true for the laws of gravity, and for the strong interactions of nuclear physics. Only the weak interactions—responsible for β -decay—do not have this symmetry. (We discussed this in some detail in Chapter 52, Vol. I.) We will for now leave out any consideration of the β -decays. Then in any physical system where β -decays are not expected to produce any appreciable effect—an example would be the emission of light by an atom—the Hamiltonian \hat{H} and the operator \hat{P} will commute. Under these circumstances we have the following proposition. If a state originally has even parity, and if you look at the physical situation at some later time, it will again have even parity. For instance, suppose an atom about to emit a photon is in a state known to have even parity. You look at the whole thing—including the photon—after the emission; it will again have even parity (likewise if you start with odd parity). This principle is called the *conservation of parity*. You can see why the words “conservation of parity” and “reflection symmetry” are closely intertwined in the quantum mechanics. Although until a few years ago it was thought that nature always conserved parity, it is now known that this is *not* true. It has been discovered to be false because the β -decay reaction does not have the inversion symmetry which is found in the other laws of physics.

Now we can prove an interesting theorem (which is true so long as we can disregard weak interactions): Any state of definite energy which is not degenerate must have a definite parity. It must have either even parity or odd parity. (Remember that we have sometimes seen systems in which several states have the same energy—we say that such states are *degenerate*. Our theorem will not apply to them.)

For a state $|\psi_0\rangle$ of definite energy, we know that

$$\hat{H}|\psi_0\rangle = E|\psi_0\rangle, \quad (17.18)$$

where E is just a number—the energy of the state. If we have *any* operator \hat{Q} which is a symmetry operator of the system we can prove that

$$\hat{Q}|\psi_0\rangle = e^{i\delta}|\psi_0\rangle \quad (17.19)$$

so long as $|\psi_0\rangle$ is a unique state of definite energy. Consider the new state $|\psi'_0\rangle$ that you get from operating with \hat{Q} . If the physics is symmetric, then $|\psi'_0\rangle$ must have the same energy as $|\psi_0\rangle$. But we have taken a situation in which there is *only one* state of that energy, namely $|\psi_0\rangle$, so $|\psi'_0\rangle$ must be the same state—it can only differ by a phase. That's the physical argument.

The same thing comes out of our mathematics. Our definition of symmetry is Eq. (17.10) or Eq. (17.11) (good for any state ψ),

$$\hat{H}\hat{Q}|\psi\rangle = \hat{Q}\hat{H}|\psi\rangle. \quad (17.20)$$

But we are considering only a state $|\psi_0\rangle$ which is a definite energy state, so that $\hat{H}|\psi_0\rangle = E|\psi_0\rangle$. Since E is just a number that floats through \hat{Q} if we want, we have

$$\hat{Q}\hat{H}|\psi_0\rangle = \hat{Q}E|\psi_0\rangle = E\hat{Q}|\psi_0\rangle.$$

So

$$\hat{H}\{\hat{Q}|\psi_0\rangle\} = E\{\hat{Q}|\psi_0\rangle\}. \quad (17.21)$$

So $|\psi'_0\rangle = \hat{Q}|\psi_0\rangle$ is also a definite energy state of \hat{H} —and with the same E . But by our hypothesis, there is only one such state; it must be that $|\psi'_0\rangle = e^{i\delta}|\psi_0\rangle$.

What we have just proved is true for any operator \hat{Q} that is a symmetry operator of the physical system. Therefore, in a situation in which we consider only electrical forces and strong interactions—and no β -decay—so that inversion symmetry is an allowed approximation, we have that $\hat{P}|\psi\rangle = e^{i\delta}|\psi\rangle$. But we have also seen that $e^{i\delta}$ must be either $+1$ or -1 . So any state of a definite energy (which is not degenerate) has got either an even parity or an odd parity.

17-3 The conservation laws

We turn now to another interesting example of an operation: a rotation. We consider the special case of an operator that rotates an atomic system by angle ϕ around the z -axis. We will call this operator† $\hat{R}_z(\phi)$. We are going to suppose that we have a physical situation where we have no influences lined up along the x - and y -axes. Any electric field or magnetic field is taken to be parallel to the z -axis‡ so that there will be no change in the *external* conditions if we rotate the whole physical system about the z -axis. For example, if we have an atom in empty space and we turn the atom around the z -axis by an angle ϕ , we have the same physical system.

Now then, there are *special states* which have the property that such an operation produces a new state which is the original state multiplied by some phase factor. Let us make a quick side remark to show you that when this is true the phase change must always be proportional to the angle ϕ . Suppose that you would rotate twice by the angle ϕ . That's the same thing as rotating by the angle 2ϕ . If a rotation by ϕ has the effect of multiplying the state $|\psi_0\rangle$ by a phase $e^{i\delta}$ so that

$$\hat{R}_z(\phi) |\psi_0\rangle = e^{i\delta} |\psi_0\rangle,$$

two such rotations in succession would multiply the state by the factor $(e^{i\delta})^2 = e^{i2\delta}$, since

$$\hat{R}_z(\phi)\hat{R}_z(\phi) |\psi_0\rangle = \hat{R}_z(\phi)e^{i\delta} |\psi_0\rangle = e^{i\delta}\hat{R}_z(\phi) |\psi_0\rangle = e^{i\delta}e^{i\delta} |\psi_0\rangle.$$

The phase change δ must be proportional to ϕ .¶ We are considering then those special states $|\psi_0\rangle$ for which

$$\hat{R}_z(\phi) |\psi_0\rangle = e^{im\phi} |\psi_0\rangle, \quad (17.22)$$

where m is some real number.

We also know the remarkable fact that *if* the system is symmetrical for a rotation around z *and if* the original state happens to have the property that (17.22) is true, then it will also have the same property later on. So this number m is a very important one. If we know its value initially, we know its value at the end of the game. It is a number which is *conserved*— m is a *constant of the motion*. The reason that we pull out m is because it hasn't anything to do with any special angle ϕ , and also because it corresponds to something in classical mechanics. In *quantum* mechanics we *choose* to call $m\hbar$ —for such states as $|\psi_0\rangle$ —the *angular momentum about the z -axis*. If we do that we find that in the limit of large systems the same quantity is equal to the z -component of the angular momentum of classical mechanics. So if we have a state for which a rotation about the z -axis just produces a phase factor $e^{im\phi}$, then we have a state of definite angular momentum about that axis—and the angular momentum is conserved. It is $m\hbar$ now and forever. Of course, you can rotate about any axis, and you get the conservation of angular momentum for the various axes. You see that the conservation of angular momentum is related to the fact that when you turn a system you get the same state with only a new phase factor.

We would like to show you how general this idea is. We will apply it to two other conservation laws which have exact correspondence in the physical ideas to the conservation of angular momentum. In classical physics we also have conservation of momentum and conservation of energy, and it is interesting to see that both of these are related in the same way to some physical symmetry.

† Very precisely, we will define $\hat{R}_z(\phi)$ as a rotation of the physical system by $-\phi$ about the z -axis, which is the same as rotating the coordinate frame by $+\phi$.

‡ We can always choose z along the direction of the field provided there is only one field at a time, and its direction doesn't change.

¶ For a fancier proof we should make this argument for small rotations ϵ . Since any angle ϕ is the sum of a suitable n number of these, $\phi = n\epsilon$, $\hat{R}_z(\phi) = [\hat{R}_z(\epsilon)]^n$ and the total phase change is n times that for the small angle ϵ , and is, therefore, proportional to ϕ .

Suppose that we have a physical system—an atom, some complicated nucleus, or a molecule, or something—and it doesn't make any difference if we take the whole system and move it over to a different place. So we have a Hamiltonian which has the property that it depends only on the *internal coordinates* in some sense, and does not depend on the *absolute position* in space. Under those circumstances there is a special symmetry operation we can perform which is a translation in space. Let's define $\hat{D}_x(a)$ as the operation of a displacement by the distance a along the x -axis. Then for any state we can make this operation and get a new state. But again there can be very special states which have the property that when you displace them by a along the x -axis you get the same state except for a phase factor. It's also possible to prove, just as we did above, that when this happens, the phase must be proportional to a . So we can write for these special states $|\psi_0\rangle$

$$\hat{D}_x(a) |\psi_0\rangle = e^{ika} |\psi_0\rangle. \quad (17.23)$$

The coefficient k , when multiplied by \hbar , is called *the x -component of the momentum*. And the reason it is called that is that this number is numerically equal to the classical momentum p_x when we have a large system. The general statement is this: If the Hamiltonian is unchanged when the system is displaced, and if the state starts with a definite momentum in the x -direction, then the momentum in the x -direction will remain the same as time goes on. The total momentum of a system before and after collisions—or after explosions or what not—will be the same.

There is another operation that is quite analogous to the displacement in space: a delay in time. Suppose that we have a physical situation where there is *nothing external* that depends on time, and we start something off at a certain moment in a given state and let it roll. Now if we were to start the same thing off again (in another experiment) two seconds later—or/say, delayed by a time τ —and if nothing in the external conditions depends on the absolute time, the development would be the same and the final state would be the same as the other final state, except that it will get there later by the time τ . Under those circumstances we can also find special states which have the property that the development in time has the special characteristic that the delayed state is just the old, multiplied by a phase factor. Once more it is clear that for these special states the phase change must be proportional to τ . We can write

$$\hat{D}_t(\tau) |\psi_0\rangle = e^{-i\omega\tau} |\psi_0\rangle. \quad (17.24)$$

It is conventional to use the negative sign in defining ω ; with this convention $\omega\hbar$ is the *energy* of the system, and it is conserved. So a system of definite energy is one which when displaced τ in time reproduces itself multiplied by $e^{-i\omega\tau}$. (That's what we have said before when we defined a quantum state of definite energy, so we're consistent with ourselves.) It means that if a system is in a state of definite energy, and if the Hamiltonian doesn't depend on t , then no matter what goes on, the system will have the same energy at all later times.

You see, therefore, the relation between the conservation laws and the symmetry of the world. Symmetry with respect to displacements in time implies the conservation of energy; symmetry with respect to position in x , y , or z implies the conservation of that component of momentum. Symmetry with respect to rotations around the x -, y -, and z -axes implies the conservation of the x -, y -, and z -components of angular momentum. Symmetry with respect to reflection implies the conservation of parity. Symmetry with respect to the interchange of two electrons implies the conservation of something we don't have a name for, and so on. Some of these principles have classical analogs and others do not. There are more conservation laws in quantum mechanics than are useful in classical mechanics—or, at least, than are usually made use of.

In order that you will be able to read other books on quantum mechanics, we must make a small technical aside—to describe the notation that people use. The operation of a displacement with respect to time is, of course, just the opera-

tion \hat{U} that we talked about before:

$$\hat{D}_t(\tau) = \hat{U}(t + \tau, t). \quad (17.25)$$

Most people like to discuss everything in terms of *infinitesimal* displacements in time, or in terms of infinitesimal displacements in space, or in terms of rotations through infinitesimal angles. Since any finite displacement or angle can be accumulated by a succession of infinitesimal displacements or angles, it is often easier to analyze first the infinitesimal case. The operator of an infinitesimal displacement Δt in time is—as we have defined it in Chapter 8—

$$\hat{D}_t(\Delta t) = 1 - \frac{i}{\hbar} \Delta t \hat{H}. \quad (17.26)$$

Then \hat{H} is analogous to the classical quantity we call energy, because if $\hat{H}|\psi\rangle$ happens to be a constant times $|\psi\rangle$ namely, $\hat{H}|\psi\rangle = E|\psi\rangle$, then that constant is the energy of the system.

The same thing is done for the other operations. If we make a small displacement in x , say by the amount Δx , a state $|\psi\rangle$ will, *in general*, go over into some other state $|\psi'\rangle$. We can write

$$|\psi'\rangle = \hat{D}_x(\Delta x)|\psi\rangle = \left(1 + \frac{i}{\hbar} \hat{p}_x \Delta x\right)|\psi\rangle, \quad (17.27)$$

since as Δx goes to zero, the $|\psi'\rangle$ should become just $|\psi\rangle$ or $\hat{D}_x(0) = 1$, and for small Δx the change of $\hat{D}_x(\Delta x)$ from 1 should be proportional to Δx . Defined this way, the operator \hat{p}_x is called the momentum operator—for the x -component, of course.

For identical reasons, people usually write for small rotations

$$\hat{R}_z(\Delta\phi)|\psi\rangle = \left(1 + \frac{i}{\hbar} \hat{J}_z \Delta\phi\right)|\psi\rangle \quad (17.28)$$

and call \hat{J}_z the operator of the z -component of angular momentum. For those special states for which $\hat{R}_z(\phi)|\psi_0\rangle = e^{im\phi}|\psi_0\rangle$, we can for any small angle—say $\Delta\phi$ —expand the right-hand side to first order in $\Delta\phi$ and get

$$\hat{R}_z(\Delta\phi)|\psi_0\rangle = e^{im\Delta\phi}|\psi_0\rangle = (1 + im\Delta\phi)|\psi_0\rangle.$$

Comparing this with the definition of \hat{J}_z in Eq. (17.28), we get that

$$\hat{J}_z|\psi_0\rangle = m\hbar|\psi_0\rangle. \quad (17.29)$$

In other words, if you operate with \hat{J}_z on a state *with a definite angular momentum* about the z -axis, you get $m\hbar$ times the same state, where $m\hbar$ is the amount of z -component of angular momentum. It is quite analogous to operating on a definite energy state with \hat{H} to get $E|\psi\rangle$.

We would now like to make some applications of the ideas of the conservation of angular momentum—to show you how they work. The point is that they are really very simple. You knew before that angular momentum is conserved. The only thing you really have to remember from this chapter is that if a state $|\psi_0\rangle$ has the property that upon a rotation through an angle ϕ about the z -axis, it becomes $e^{im\phi}|\psi_0\rangle$; it has a z -component of angular momentum equal to $m\hbar$. That's all we will need to do a number of interesting things.

17-4 Polarized light

First of all we would like to check on one idea. In Section 11-4 we showed that when RHC polarized light is viewed in a frame rotated by the angle ϕ about the z -axis† it gets multiplied by $e^{i\phi}$. Does that mean then that the photons of light

† Sorry! This angle is the negative of the one we used in Section 11-4.

that are right circularly polarized carry an angular momentum of *one* unit† along the *z*-axis? *Indeed it does.* It also means that if we have a beam of light containing a large number of photons all circularly polarized the same way—as we would have in a classical beam—it will carry angular momentum. If the total energy carried by the beam in a certain time is W , then there are $N = W/\hbar\omega$ photons. Each one carries the angular momentum \hbar , so there is a total angular momentum of

$$J_z = N\hbar = \frac{W}{\omega}. \quad (17.30)$$

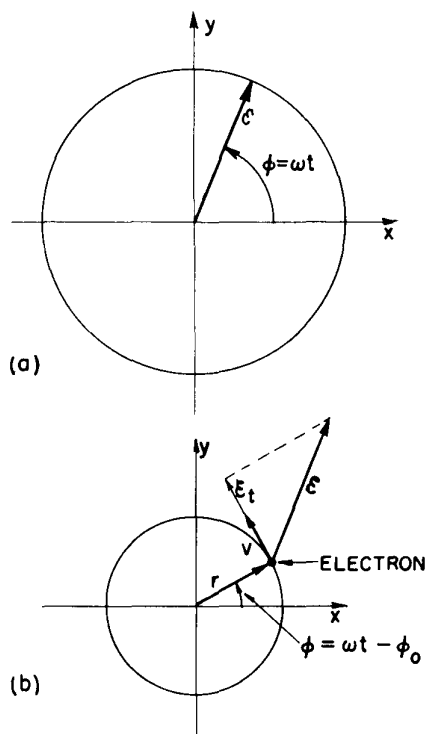


Fig. 17-5. (a) The electric field \mathcal{E} in a circularly polarized light wave. (b) The motion of an electron being driven by the circularly polarized light.

Can we prove classically that light which is right circularly polarized carries an energy and angular momentum in proportion to W/ω ? That should be a classical proposition if everything is right. Here we have a case where we can go from the quantum thing to the classical thing. We should see if the classical physics checks. It will give us an idea whether we have a right to call m the angular momentum. Remember what right circularly polarized light is, classically. It's described by an electric field with an oscillating *x*-component and an oscillating *y*-component 90° out of phase so that the resultant electric vector \mathcal{E} goes in a circle—as drawn in Fig. 17-5(a). Now suppose that such light shines on a wall which is going to absorb it—or at least some of it—and consider an atom in the wall according to the classical physics. We have often described the motion of the electron in the atom as a harmonic oscillator which can be driven into oscillation by an external electric field. We'll suppose that the atom is isotropic, so that it can oscillate equally well in the *x*- or *y*-directions. Then in the circularly polarized light, the *x*-displacement and the *y*-displacement are the same, but one is 90° behind the other. The net result is that the electron moves in a circle, as shown in Fig. 17-5(b). The electron is displaced at some displacement r from its equilibrium position at the origin and goes around with some phase lag with respect to the vector \mathcal{E} . The relation between \mathcal{E} and r might be as shown in Fig. 17-5(b). As time goes on, the electric field rotates and the displacement rotates with the same frequency, so their relative orientation stays the same. Now let's look at the work being done on this electron. The rate that energy is being put into this electron is v , its velocity, times the component of $q\mathcal{E}$ parallel to the velocity:

$$\frac{dW}{dt} = q\mathcal{E}_t v. \quad (17.31)$$

But look, there is angular momentum being poured into this electron, because there is always a torque about the origin. The torque is $q\mathcal{E}_t r$, which must be equal to the rate of change of angular momentum dJ_z/dt :

$$\frac{dJ_z}{dt} = q\mathcal{E}_t r. \quad (17.32)$$

Remembering that $v = \omega r$, we have that

$$\frac{dJ_z}{dW} = \frac{1}{\omega}.$$

Therefore, if we integrate the total angular momentum which is absorbed, it is proportional to the total energy—the constant of proportionality being $1/\omega$, which agrees with Eq. (17.30). Light does carry angular momentum—1 unit (times \hbar) if it is right circularly polarized along the *z*-axis, and -1 unit along the *z*-axis if it is left circularly polarized.

Now let's ask the following question: If light is linearly polarized in the *x*-direction, what is its angular momentum? Light polarized in the *x*-direction can be represented as the superposition of RHC and LHC polarized light. Therefore, there is a certain amplitude that the angular momentum is $+\hbar$ and another

† It is usually very convenient to measure angular momentum of atomic systems in units of \hbar . Then you can say that a spin one-half particle has angular momentum $\pm 1/2$ with respect to any axis. Or, in general, that the *z*-component of angular momentum is m . You don't need to repeat the \hbar all the time.

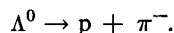
amplitude that the angular momentum is $-\hbar$, so it doesn't have a *definite* angular momentum. It has an amplitude to appear with $+\hbar$ and an equal amplitude to appear with $-\hbar$. The interference of these two amplitudes produces the linear polarization, but it has *equal* probabilities to appear with plus or minus one unit of angular momentum. Macroscopic measurements made on a beam of linearly polarized light will show that it carries zero angular momentum, because in a large number of photons there are nearly equal numbers of RHC and LHC photons contributing opposite amounts of angular momentum—the average angular momentum is zero. And in the classical theory you don't find the angular momentum unless there is some circular polarization.

We have said that any spin-one particle can have three values of J_z , namely $+1, 0, -1$ (the three states we saw in the Stern-Gerlach experiment). But light is screwy; it has only two states. It does not have the zero case. This strange lack is related to the fact that light cannot stand still. For a particle of spin j which is standing still, there must be the $2j + 1$ possible states with values of j_z going in steps of 1 from $-j$ to $+j$. But it turns out that for something of spin j with zero mass only the states with the components $+j$ and $-j$ along the direction of motion exist. For example, light does not have three states, but only two—although a photon is still an object of spin one. How is this consistent with our earlier proofs—based on what happens under rotations in space—that for spin-one particles three states are necessary? For a particle at rest, rotations can be made about any axis without changing the momentum state. Particles with zero rest mass (like photons and neutrinos) cannot be at rest; only rotations about the axis along the direction of motion do not change the momentum state. Arguments about rotations around one axis only are insufficient to prove that three states are required, given that one of them varies as $e^{i\phi}$ under rotations by the angle ϕ .†

One further side remark. For a zero rest mass particle, in general, *only one* of the two spin states with respect to the line of motion ($+j, -j$) is really necessary. For neutrinos—which are spin one-half particles—only the states with the component of angular momentum *opposite* to the direction of motion ($-\hbar/2$) exist in nature [and only *along* the motion ($+\hbar/2$) for antineutrinos]. When a system has inversion symmetry (so that parity is conserved, as it is for light) both components ($+j$, and $-j$) are required.

17-5 The disintegration of the Λ^0

Now we want to give an example of how we use the theorem of conservation of angular momentum in a specifically quantum physical problem. We look at break-up of the lambda particle (Λ^0), which disintegrates into a proton and a π^- meson by a “weak” interaction:



Assume we know that the pion has spin zero, that the proton has spin one-half, and that the Λ^0 has spin one-half. We would like to solve the following problem: Suppose that a Λ^0 were to be produced in a way that caused it to be completely polarized—by which we mean that its spin is, say “up,” with respect to some suitably chosen z -axis—see Fig. 17-6(a). The question is, with what probability will it disintegrate so that the proton goes off at an angle θ with respect to the z -axis—as in Fig. 17-6(b)? In other words, what is the angular distribution of the disintegrations? We will look at the disintegration in the coordinate system in which the Λ^0 is at rest—we will measure the angles in this rest frame; then they can always be transformed to another frame if we want.

† We have tried to find at least a proof that the component of angular momentum along the direction of motion must for a zero mass particle be an integral multiple of $\hbar/2$ —and not something like $\hbar/3$. Even using all sorts of properties of the Lorentz transformation and what not, we failed. Maybe it's not true. We'll have to talk about it with Prof. Wigner, who knows all about such things.

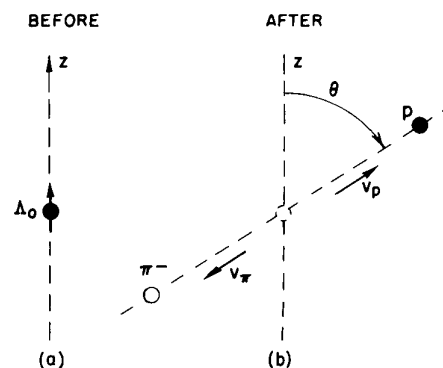


Fig. 17-6. A Λ^0 with spin “up” decays into a proton and a pion (in the CM system). What is the probability that the proton will go off at the angle θ ?

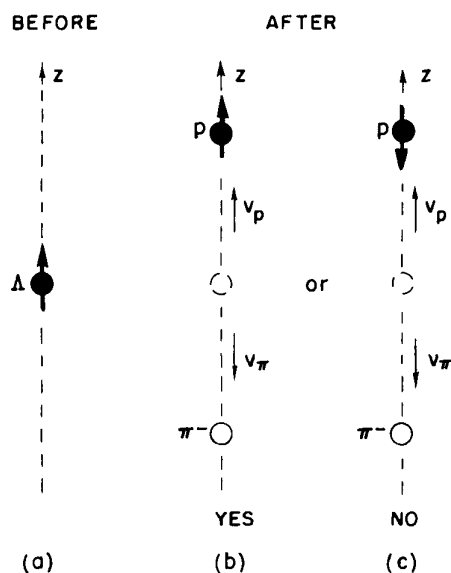


Fig. 17-7. Two possibilities for the decay of a spin "up" Λ^0 with the proton going along the $+z$ -axis. Only (b) conserves angular momentum.

We begin by looking at the special circumstance in which the proton is emitted into a small solid angle $\Delta\Omega$ along the z -axis (Fig. 17-7). Before the disintegration we have a Λ^0 with its spin "up," as in part (a) of the figure. After a short time—for reasons unknown to this day, except that they are connected with the weak decays—the Λ^0 explodes into a proton and a pion. Suppose the proton goes up along the $+z$ -axis. Then, from the conservation of momentum, the pion must go down. Since the proton is a spin one-half particle, its spin must be either "up" or "down"—there are, in principle, the two possibilities shown in parts (b) and (c) of the figure. The conservation of angular momentum, however, requires that the proton have spin "up." This is most easily seen from the following argument. A particle moving along the z -axis cannot contribute any angular momentum about this axis by virtue of its motion; therefore, only the spins can contribute to J_z . The spin angular momentum about the z -axis is $+\hbar/2$ before the disintegration, so it must also be $+\hbar/2$ afterward. We can say that since the pion has no spin, the proton spin must be "up."

If you are worried that arguments of this kind may not be valid in quantum mechanics, we can take a moment to show you that they are. The initial state (before the disintegration), which we can call $|\Lambda^0, \text{spin } +z\rangle$ has the property that if it is rotated about the z -axis by the angle ϕ , the state vector gets multiplied by the phase factor $e^{+i\phi/2}$. (In the rotated system the state vector is $e^{+i\phi/2} |\Lambda^0, \text{spin } +z\rangle$.) That's what we mean by spin "up" for a spin one-half particle. Since nature's behavior doesn't depend on our choice of axes, the final state (the proton plus pion) must have the same property. We could write the final state as, say,

$$|\text{proton going } +z, \text{spin } +z; \text{pion going } -z\rangle.$$

But we really do not need to specify the pion motion, since in the frame we have chosen the pion always moves opposite the proton; we can simplify our description of the final state to

$$|\text{proton going } +z, \text{spin } +z\rangle.$$

Now what happens to this state vector if we rotate the coordinates about the z -axis by the angle ϕ ?

Since the proton and pion are moving along the z -axis, their motion isn't changed by the rotation. (That's why we picked this special case; we couldn't make the argument otherwise.) Also, nothing happens to the pion, because it is spin zero. The proton, however, has spin one-half. If its spin is "up" it will contribute a phase change of $e^{+i\phi/2}$ in response to the rotation. (If its spin were "down" the phase change due to the proton would be $e^{-i\phi/2}$.) But the phase change with rotation before and after the excitement must be the same if angular momentum is to be conserved. (And it will be, since there are no outside influences in the Hamiltonian.) So the only possibility is that the proton spin will be "up." If the proton goes up, its spin must also be "up."

We conclude, then, that the conservation of angular momentum permits the process shown in part (b) of Fig. 17-7, but does not permit the process shown in part (c). Since we know that the disintegration occurs, there is some amplitude for process (b)—proton going up with spin "up." We'll let a stand for the amplitude that the disintegration occurs in this way in any infinitesimal interval of time.[†]

Now let's see what would happen if the Λ^0 spin were initially "down." Again we ask about the decays in which the proton goes up along the z -axis, as shown in Fig. 17-8. You will appreciate that in this case the proton must have spin "down" if angular momentum is conserved. Let's say that the amplitude for such a disintegration is b .

We can't say anything more about the two amplitudes a and b . They depend on the inner machinery of Λ^0 , and the weak decays, and nobody yet knows how to

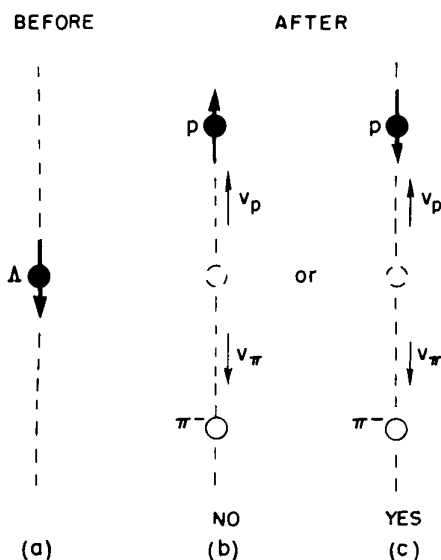


Fig. 17-8. The decay along the z -axis for a Λ^0 with spin "down."

[†] We are now assuming that the machinery of the quantum mechanics is sufficiently familiar to you that we can speak about things in a physical way without taking the time to write down all the mathematical details. In case what we are saying here is not clear to you, we have put some of the missing details in a note at the end of the section.

calculate them. We'll have to get them from experiment. But with just these two amplitudes we *can* find out all we want to know about the angular distribution of the disintegration. We only have to be careful always to define completely the states we are talking about.

We want to know the probability that the proton will go off at the angle θ with respect to the z -axis (into a small solid angle $\Delta\Omega$) as drawn in Fig. 17-6. Let's put a new z -axis in this direction and call it the z' -axis. We know how to analyze what happens along this axis. With respect to this new axis, the Λ^0 no longer has its spin "up," but has a certain amplitude to have its spin "up" and another amplitude to have its spin "down." We have already worked these out in Chapter 6, and again in Chapter 10, Eq. (10.30). The amplitude to be spin "up" is $\cos \theta/2$, and the amplitude to be spin "down" is† $-\sin \theta/2$. When the Λ^0 spin is "up" along the z' -axis it will emit a proton in the $+z'$ -direction with the amplitude a . So the amplitude to find an "up"-spinning proton coming out along the z' -direction is

$$a \cos \frac{\theta}{2}. \quad (17.33)$$

Similarly, the amplitude to find a "down"-spinning proton coming along the positive z' -axis is

$$-b \sin \frac{\theta}{2}. \quad (17.34)$$

The two processes that these amplitudes refer to are shown in Fig. 17-9.

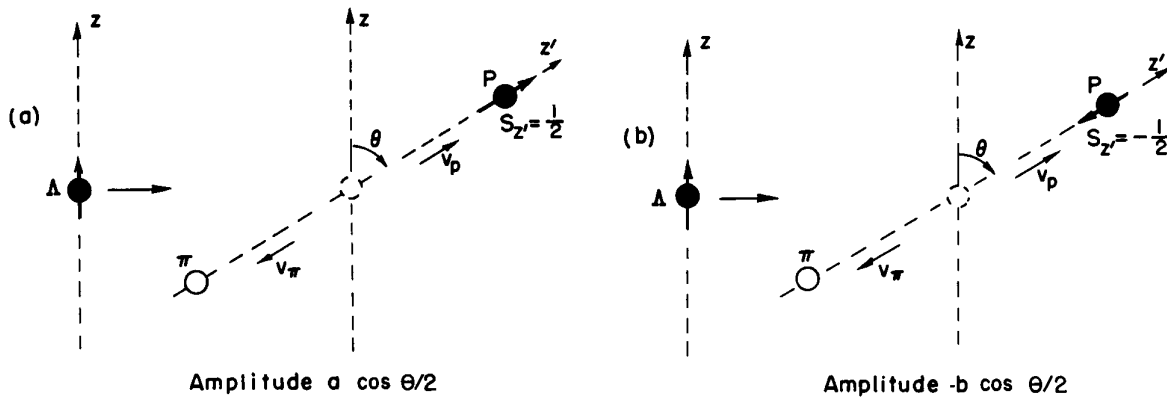


Fig. 17-9. Two possible decay states for the Λ^0 .

Let's now ask the following easy question. If the Λ^0 has spin up along the z -axis, what is the probability that the decay proton will go off at the angle θ ? The two spin states ("up" or "down" along z') are distinguishable even though we are not going to look at them. So to get the probability we square the amplitudes and add. The probability $f(\theta)$ of finding a proton in a small solid angle $\Delta\Omega$ at θ is

$$f(\theta) = |a|^2 \cos^2 \frac{\theta}{2} + |b|^2 \sin^2 \frac{\theta}{2}. \quad (17.35)$$

Remembering that $\sin^2 \theta/2 = \frac{1}{2}(1 - \cos \theta)$ and that $\cos^2 \theta/2 = \frac{1}{2}(1 + \cos \theta)$, we can write $f(\theta)$ as

$$f(\theta) = \left(\frac{|a|^2 + |b|^2}{2} \right) + \left(\frac{|a|^2 - |b|^2}{2} \right) \cos \theta. \quad (17.36)$$

† We have chosen to let z' be in the xz -plane and use the matrix elements for $R_y(\theta)$. You would get the same answer for any other choice.

The angular distribution has the form

$$f(\theta) = \beta(1 + \alpha \cos \theta). \quad (17.37)$$

The probability has one part that is independent of θ and one part that varies linearly with $\cos \theta$. From measuring the angular distribution we can get α and β , and therefore, $|a|$ and $|b|$.

Now there are many other questions we can answer. Are we interested only in protons with spin “up” along the *old* z -axis? Each of the terms in (17-33) and (17-34) will give an amplitude to find a proton with spin “up” and with spin “down” with respect to the z' -axis ($+z'$ and $-z'$). Spin “up” with respect to the old axis $|+z\rangle$ can be expressed in terms of the base states $|+z'\rangle$ and $|-z'\rangle$. We can then combine the two amplitudes (17.33) and (17.34) with the proper coefficients ($\cos \theta/2$ and $-\sin \theta/2$) to get the total amplitude

$$\left(a \cos^2 \frac{\theta}{2} + b \sin^2 \frac{\theta}{2} \right).$$

Its square is the probability that the proton comes out at the angle θ with its spin the same as the Λ^0 (“up” along the z -axis).

If parity were conserved, we could say one more thing. The disintegration of Fig. 17-8 is just the reflection—in say, the y z -plane of the disintegration of Fig. 17-7.† If parity were conserved, b would have to be equal to a or to $-a$. Then the coefficient α of (17.37) would be zero, and the disintegration would be equally likely to occur in all directions.

The experimental results show, however, that there *is* an asymmetry in the disintegration. The measured angular distribution does go as $\cos \theta$ as we predict—and not as $\cos^2 \theta$ or any other power. In fact, since the angular distribution has this form, we can deduce from these measurements that the spin of the Λ^0 is $1/2$. Also, we see that parity is not conserved. In fact, the coefficient α is found experimentally to be -0.62 ± 0.05 , so b is about twice as large as a . The lack of symmetry under a reflection is quite clear.

You see how much we can get from the conservation of angular momentum. We will give some more examples in the next chapter.

Parenthetical note. By the amplitude a in this section we mean the amplitude that the state $|\text{proton going } +z, \text{ spin } +z\rangle$ is generated in an infinitesimal time dt from the state $|\Lambda, \text{ spin } +z\rangle$, or, in other words, that

$$\langle \text{proton going } +z, \text{ spin } +z | H | \Lambda, \text{ spin } +z \rangle = i\hbar a, \quad (17.38)$$

where H is the Hamiltonian of the world—or, at least, of whatever is responsible for the Λ -decay. The conservation of angular momentum means that the Hamiltonian must have the property that

$$\langle \text{proton going } +z, \text{ spin } -z | H | \Lambda, \text{ spin } +z \rangle = 0. \quad (17.39)$$

By the amplitude b we mean that

$$\langle \text{proton going } +z, \text{ spin } -z | H | \Lambda, \text{ spin } -z \rangle = i\hbar b. \quad (17.40)$$

Conservation of angular momentum implies that

$$\langle \text{proton going } +z, \text{ spin } +z | H | \Lambda, \text{ spin } -z \rangle = 0. \quad (17.41)$$

If the amplitudes written in (17.33) and (17.34) are not clear, we can express them more mathematically as follows. By (17.33) we intend the amplitude that the Λ with spin along $+z$ will disintegrate into a proton moving along the $+z'$ -direction with its spin also in the $+z'$ -direction, namely the amplitude

$$\langle \text{proton going } +z', \text{ spin } +z' | H | \Lambda, \text{ spin } +z \rangle. \quad (17.42)$$

By the general theorems of quantum mechanics, this amplitude can be written as

$$\sum \langle \text{proton going } +z', \text{ spin } +z' | H | \Lambda, i \rangle \langle \Lambda, i | \Lambda, \text{ spin } +z \rangle, \quad (17.43)$$

† Remembering that the spin is an axial vector and flips over in the reflection.

where the sum is to be taken over the base states $|\Lambda, i\rangle$ of the Λ -particle at rest. Since the Λ -particle is spin one-half, there are two such base states which can be in any reference base we wish. If we use for base states spin “up” and spin “down” *with respect to z'* ($+z'$, $-z'$), the amplitude of (17.43) is equal to the sum

$$\begin{aligned} &\langle \text{proton going } +z', \text{ spin } +z' | H | \Lambda, +z' \rangle \langle \Lambda, +z' | \Lambda, +z \rangle \\ &+ \langle \text{proton going } +z', \text{ spin } +z' | H | \Lambda, -z' \rangle \langle \Lambda, -z' | \Lambda, +z \rangle. \end{aligned} \quad (17.44)$$

The first factor of the first term is a , and the first factor of the second term is zero—from the definition of (17.38), and from (17.41), which in turn follows from angular momentum conservation. The remaining factor $\langle \Lambda, +z' | \Lambda, +z \rangle$ of the first term is just the amplitude that a spin one-half particle which has spin “up” along one axis will also have spin “up” along an axis tilted at the angle θ , which is $\cos \theta/2$ —see Table 6-2. So (17.44) is just $a \cos \theta/2$, as we wrote in (17.33). The amplitude of (17.34) follows from the same kind of arguments for a spin “down” Λ -particle.

17-6 Summary of the rotation matrices

We would like now to bring together in one place the various things we have learned about the rotations for particles of spin one-half and spin one—so they will be convenient for future reference. On the next page you will find tables of the two rotation matrices $R_z(\phi)$ and $R_y(\theta)$ for spin one-half particles, for spin-one particles, and for photons (spin-one particles with zero rest mass). For each spin we will give the terms of the matrix $\langle j | R | i \rangle$ for rotations about the z -axis or the y -axis. They are, of course, exactly equivalent to the amplitudes like $\langle +T | 0 S \rangle$ we have used in earlier chapters. We mean by $R_z(\phi)$ that the state is projected into a new coordinate system which is rotated through the angle ϕ about the z -axis—using always the right-hand rule to define the positive sense of the rotation. By $R_y(\theta)$ we mean that the reference axes are rotated by the angle θ about the y -axis. Knowing these two rotations, you can, of course, work out any arbitrary rotation. As usual, we write the matrix elements so that the state on the *left* is a base state of the *new* (rotated) frame and the state on the *right* is a base state of the old (unrotated) frame. You can interpret the entries in the tables in many ways. For instance, the entry $e^{-i\phi/2}$ in Table 17-1 means that the matrix element $\langle - | R | - \rangle = e^{-i\phi/2}$. It also means that $\hat{R} | - \rangle = e^{-i\phi/2} | - \rangle$, or that $\langle - | \hat{R} = \langle - | e^{-i\phi/2}$. It's all the same thing.

Table 17-1

Rotation matrices for spin one-half

Two states: $|+\rangle$, "up" along the z-axis, $m = +1/2$
 $|-\rangle$, "down" along the z-axis, $m = -1/2$

$R_z(\phi)$	$ +\rangle$	$ -\rangle$
$\langle+ $	$e^{+i\phi/2}$	0
$\langle- $	0	$e^{-i\phi/2}$

$R_y(\theta)$	$ +\rangle$	$ -\rangle$
$\langle+ $	$\cos \theta/2$	$\sin \theta/2$
$\langle- $	$-\sin \theta/2$	$\cos \theta/2$

Table 17-2

Rotation matrices for spin one

Three states: $|+\rangle$, $m = +1$
 $|0\rangle$, $m = 0$
 $|-\rangle$, $m = -1$

$R_z(\phi)$	$ +\rangle$	$ 0\rangle$	$ -\rangle$
$\langle+ $	$e^{+i\phi}$	0	0
$\langle 0 $	0	1	0
$\langle- $	0	0	$e^{-i\phi}$

$R_y(\theta)$	$ +\rangle$	$ 0\rangle$	$ -\rangle$
$\langle+ $	$\frac{1}{2}(1 + \cos \theta)$	$+\frac{1}{\sqrt{2}} \sin \theta$	$\frac{1}{2}(1 - \cos \theta)$
$\langle 0 $	$-\frac{1}{\sqrt{2}} \sin \theta$	$\cos \theta$	$+\frac{1}{\sqrt{2}} \sin \theta$
$\langle- $	$\frac{1}{2}(1 - \cos \theta)$	$-\frac{1}{\sqrt{2}} \sin \theta$	$\frac{1}{2}(1 + \cos \theta)$

Table 17-3

Photons

Two states: $|R\rangle = \frac{1}{\sqrt{2}}(|x\rangle + i|y\rangle)$, $m = +1$ (RHC polarized)

$|L\rangle = \frac{1}{\sqrt{2}}(|x\rangle - i|y\rangle)$, $m = -1$ (LHC polarized)

$R_z(\phi)$	$ R\rangle$	$ L\rangle$
$\langle R $	$e^{+i\phi}$	0
$\langle L $	0	$e^{-i\phi}$

Angular Momentum

18-1 Electric dipole radiation

In the last chapter we developed the idea of the conservation of angular momentum in quantum mechanics, and showed how it might be used to predict the angular distribution of the photon from the disintegration of the Λ -particle. We want now to give you a number of other, similar, illustrations of the consequences of momentum conservation in atomic systems. Our first example is the radiation of light from an atom. The conservation of angular momentum (among other things) will determine the polarization and angular distribution of the emitted photons.

Suppose we have an atom which is in an excited state of definite angular momentum—say with a spin of one—and it makes a transition to a state of angular momentum zero at a lower energy, emitting a photon. The problem is to figure out the angular distribution and polarization of the photons. (This problem is almost exactly the same as the Λ^0 disintegration, except that we have spin-one instead of spin one-half particles.) Since the upper state of the atom is spin one, there are three possibilities for its z -component of angular momentum. The value of m could be $+1$, or 0 , or -1 . We will take $m = +1$ for our example. Once you see how it goes, you can work out the other cases. We suppose that the atom is sitting with its angular momentum along the $+z$ -axis—as in Fig. 18-1(a)—and ask with what amplitude it will emit right circularly polarized light upward along the z -axis, so that the atom ends up with zero angular momentum—as shown in part (b) of the figure. Well, we don't know the answer to that. But we do know that right circularly polarized light has one unit of angular momentum about its direction of propagation. So after the photon is emitted, the situation would have to be as shown in Fig. 18-1(b)—the atom is left with zero angular momentum

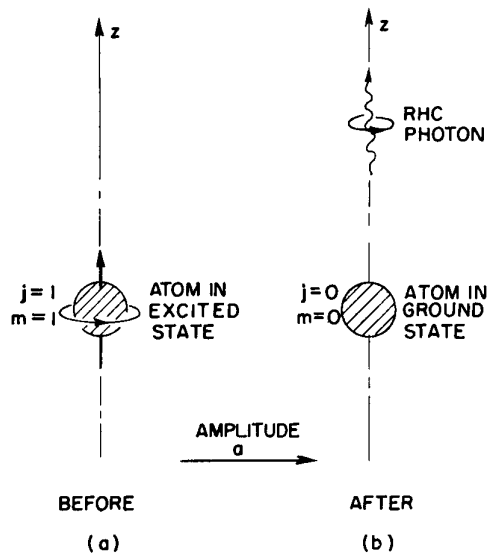


Fig. 18-1. An atom with $m = +1$ emits a RHC photon along the $+z$ -axis.

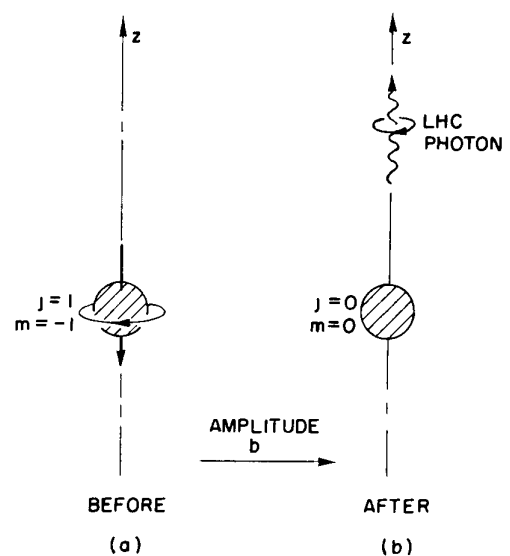


Fig. 18-2. An atom with $m = -1$ emits a LHC photon along the $+z$ -axis.

18-1 Electric dipole radiation

18-2 Light scattering

18-3 The annihilation of positronium

18-4 Rotation matrix for any spin

18-5 Measuring a nuclear spin

18-6 Composition of angular momentum

Added Note 1: Derivation of the rotation matrix

Added Note 2: Conservation of parity in photon emission

about the z -axis, since we have assumed an atom whose lower state is spin zero. We will let a stand for the amplitude for such an event. More precisely, we let a be the amplitude to emit a photon into a certain small solid angle $\Delta\Omega$, centered on the z -axis, during a time dt . Notice that the amplitude to emit a LHC photon in the same direction is zero. The net angular momentum about the z -axis would be -1 for such a photon and zero for the atom for a total of -1 , which would not conserve angular momentum.

Similarly, if the spin of the atom is initially “down” (-1 along the z -axis), it can emit only a LHC polarized photon in the direction of the $+z$ -axis, as shown in Fig. 18-2. We will let b stand for the amplitude for this event—meaning again the amplitude that the photon goes into a certain solid angle $\Delta\Omega$. On the other hand, if the atom is in the $m = 0$ state, it cannot emit a photon in the $+z$ -direction at all, because a photon can have only the angular momentum $+1$ or -1 along its direction of motion.

Next, we can show that b is related to a . Suppose we perform an inversion of the situation in Fig. 18-1, which means that we should imagine what the system would look like if we were to move each part of the system to an equivalent point on the opposite side of the origin. This does *not* mean that we should reflect the angular momentum vectors, because they are artificial. We should, rather, invert the actual character of the motion that would correspond to such an angular momentum. In Fig. 18-3(a) and (b) we show what the process of Fig. 18-1 looks like before and after an inversion with respect to the center of the atom. Notice that the sense of rotation of the atom is unchanged.[†] In the inverted system of Fig. 18-3(b) we have an atom with $m = +1$ emitting a LHC photon downward.

If we now rotate the system of Fig. 18-3(b) by 180° about the x - or y -axis, it becomes identical to Fig. 18-2. The combination of the inversion and rotation turns the second process into the first. Using Table 17-2, we see that a rotation of 180° about the y -axis just throws an $m = -1$ state into an $m = +1$ state, so the amplitude b must be equal to the amplitude a *except for a possible sign change due to the inversion*. The sign change in the inversion will depend on the parities of the initial and final state of the atom.

In atomic processes, parity is conserved, so the parity of the whole system must be the same before and after the photon emission. What happens will depend on whether the parities of the initial and final states of the atom are even or odd—the angular distribution of the radiation will be different for different cases. We will take the common case of *odd* parity for the initial state and *even* parity for the final state; it will give what is called “electric dipole radiation.” (If the initial and final states have the same parity we say there is “magnetic dipole radiation,” which has the character of the radiation from an oscillating current in a loop.) If the parity of the initial state is odd, its amplitude reverses its sign in the inversion which takes the system from (a) to (b) of Fig. 18-3. The final state of the atom has even parity, so its amplitude doesn’t change sign. If the reaction is going to conserve parity, the amplitude b must be equal to a in magnitude but of the opposite sign.

We conclude that if the amplitude is a that an $m = +1$ state will emit a photon upward, then for the assumed parities of the initial and final states the amplitude that an $m = -1$ state will emit a LHC photon upward is $-a$.[‡]

We have all we need to know to find the amplitude for a photon to be emitted at any angle θ with respect to the z -axis. Suppose we have an atom originally polarized with $m = +1$. We can resolve this state into $+1$, 0 , and -1 states with respect to a new z' -axis in the direction of the photon emission. The amplitudes for these three states are just the ones given in the lower half of Table 17-2.

[†] When we change x, y, z into $-x, -y, -z$, you might think that all vectors get reversed. That is true for *polar* vectors like displacements and velocities, but *not* for an *axial* vector like angular momentum—or any vector which is derived from a cross product of two polar vectors. Axial vectors have the same components after an inversion.

[‡] Some of you may object to the argument we have just made, on the basis that the final states we have been considering do not have a definite parity. You will find in Added Note 2 at the end of this chapter another demonstration, which you may prefer.

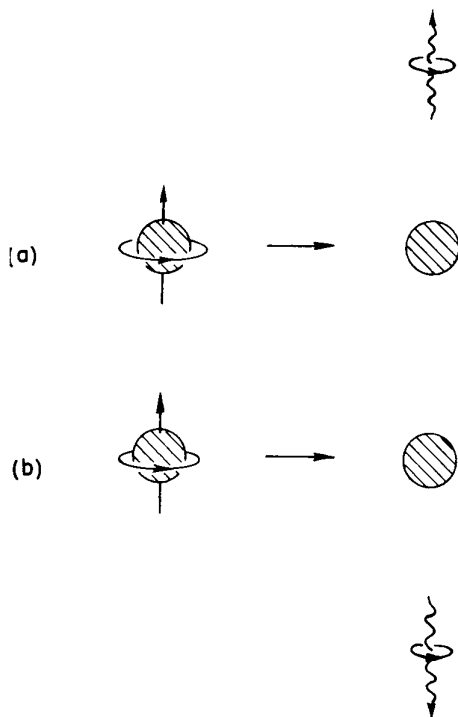


Fig. 18-3. If the process of (a) is transformed by an inversion through the center of the atom, it appears as in (b).

The amplitude that a RHC photon is emitted in the direction θ is then a times the amplitude to have $m = +1$ in that direction, namely,

$$a \langle + | R_y(\theta) | + \rangle = \frac{a}{2} (1 + \cos \theta). \quad (18.1)$$

The amplitude that a LHC photon is emitted in the same direction is $-a$ times the amplitude to have $m = -1$ in the new direction. Using Table 17-2, it is

$$-a \langle - | R_y(\theta) | + \rangle = \frac{a}{-2} (1 - \cos \theta). \quad (18.2)$$

If you are interested in other polarizations you can find out the amplitude for them from the superposition of these two amplitudes. To get the intensity of any component as a function of angle, you must, of course, take the absolute square of the amplitudes.

18-2 Light scattering

Let's use these results to solve a somewhat more complicated problem—but also one which is somewhat more real. We suppose that the same atoms are sitting in their ground state ($j = 0$), and *scatter* an incoming beam of light. Let's say that the light is going initially in the $+z$ -direction, so that we have photons coming up to the atom *from* the $-z$ -direction, as shown in Fig. 18-4(a). We can consider the scattering of light as a two-step process: The photon is absorbed, and then is re-emitted. If we start with a RHC photon as in Fig. 18-4(a), and angular momentum is conserved, the atom will be in an $m = +1$ state after the absorption—as shown in Fig. 18-4(b). We call the amplitude for this process c . The atom can then emit a RHC photon in the direction θ —as in Fig. 18-4(c). The total amplitude that a RHC photon is scattered in the direction θ is just c times (18.1). Let's call this scattering amplitude $\langle R' | S | R \rangle$; we have

$$\langle R' | S | R \rangle = \frac{ac}{2} (1 + \cos \theta). \quad (18.3)$$

There is also an amplitude that a RHC photon will be absorbed and that a LHC photon will be emitted. The product of the two amplitudes is the amplitude $\langle L' | S | R \rangle$ that a RHC photon is scattered as a LHC photon. Using (18.2), we have

$$\langle L' | S | R \rangle = -\frac{ac}{2} (1 - \cos \theta). \quad (18.4)$$

Now let's ask about what happens if a LHC photon comes in. When it is absorbed, the atom will go into an $m = -1$ state. By the same kind of arguments we used in the preceding section, we can show that this amplitude must be $-c$. The amplitude that an atom in the $m = -1$ state will emit a RHC photon at the angle θ is a times the amplitude $\langle + | R_y(\theta) | - \rangle$, which is $\frac{1}{2}(1 - \cos \theta)$. So we have

$$\langle R' | S | L \rangle = -\frac{ac}{2} (1 - \cos \theta). \quad (18.5)$$

Finally, the amplitude for a LHC photon to be scattered as a LHC photon is

$$\langle L' | S | L \rangle = \frac{ac}{2} (1 + \cos \theta). \quad (18.6)$$

(There are two minus signs which cancel.)

If we make a measurement of the scattered *intensity* for any given combination of circular polarizations it will be proportional to the square of one of our four amplitudes. For instance, with an incoming beam of RHC light the intensity of the RHC light in the scattered radiation will vary as $(1 + \cos \theta)^2$.

That's all very well, but suppose we start out with *linearly* polarized light. What then? If we have x -polarized light, it can be represented as a superposition

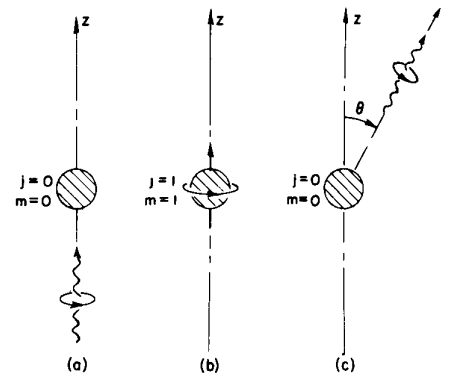


Fig. 18-4. The scattering of light by an atom seen as a two-step process.

of RHC and LHC light. We write (see Section 11-4)

$$|x\rangle = \frac{1}{\sqrt{2}} (|R\rangle + |L\rangle). \quad (18.7)$$

Or, if we have y -polarized light, we would have

$$|y\rangle = -\frac{i}{\sqrt{2}} (|R\rangle - |L\rangle). \quad (18.8)$$

Now what do you want to know? Do you want the amplitude that an x -polarized photon will scatter into a RHC photon at the angle θ ? You can get it by the usual rule for combining amplitudes. First, multiply (18.7) by $\langle R' | S$ to get

$$\langle R' | S | x \rangle = \frac{1}{\sqrt{2}} (\langle R' | S | R \rangle + \langle R' | S | L \rangle), \quad (18.9)$$

and then use (18.3) and (18.5) for the two amplitudes. You get

$$\langle R' | S | x \rangle = \frac{ac}{\sqrt{2}} \cos \theta. \quad (18.10)$$

If you wanted the amplitude that an x -photon would scatter into a LHC photon, you would get

$$\langle L' | S | x \rangle = \frac{ac}{\sqrt{2}} \cos \theta. \quad (18.11)$$

Finally, suppose you wanted to know the amplitude that an x -polarized photon will scatter while keeping its x -polarization. What you want is $\langle x' | S | x \rangle$. This can be written as

$$\langle x' | S | x \rangle = \langle x' | R' \rangle \langle R' | S | x \rangle + \langle x' | L' \rangle \langle L' | S | x \rangle. \quad (18.12)$$

If you then use the relations

$$|R'\rangle = \frac{1}{\sqrt{2}} (|x'\rangle + i|y'\rangle), \quad (18.13)$$

$$|L'\rangle = \frac{1}{\sqrt{2}} (|x'\rangle - i|y'\rangle), \quad (18.14)$$

it follows that

$$\langle x' | R' \rangle = \frac{1}{\sqrt{2}}, \quad (18.15)$$

$$\langle x' | L' \rangle = \frac{1}{\sqrt{2}}. \quad (18.16)$$

So you get that

$$\langle x' | S | x \rangle = ac \cos \theta. \quad (18.17)$$

The answer is that a beam of x -polarized light will be scattered at the direction θ (in the xz -plane) with an *intensity* proportional to $\cos^2 \theta$. If you ask about y -polarized light, you find that

$$\langle y' | S | x \rangle = 0. \quad (18.18)$$

So the scattered light is completely polarized in the x -direction.

Now we notice something interesting. The results (18.17) and (18.18) correspond exactly to the classical theory of light scattering we gave in Vol. 1, Section 32-6, where we imagined that the electron was bound to the atom by a linear restoring force—so that it acted like a classical oscillator. Perhaps you are thinking: “It’s so much easier in the classical theory; if it gives the right answer why bother with the quantum theory?” For one thing, we have considered so far only the special—though common—case of an atom with a $j = 1$ excited state and a $j = 0$ ground state. If the excited state had spin two, you would get a different result. Also, there is no reason why the model of an electron attached to a

spring and driven by an oscillating electric field should work for a single photon. But we have found that it does in fact work, and that the polarization and intensities come out right. So in a certain sense we are bringing the whole course around to the real truth. Whereas we have, in Vol. I, done the theory of the index of refraction, and of light scattering, by the classical theory, we have now shown that the quantum theory gives the same result for the most common case. In effect we have now done the polarization of sky light, for instance, by quantum mechanical arguments, which is the only truly legitimate way.

It should be, of course, that all the classical theories which work are supported ultimately by legitimate quantum arguments. Naturally, those things which we have spent a great deal of time in explaining to you were selected from just those parts of classical physics which still maintain validity in quantum mechanics. You'll notice that we did not discuss in great detail any model of the atom which has electrons going around in orbits. That's because such a model doesn't give results which agree with the quantum mechanics. But the electron on a spring—which is not, in a sense, at all the way an atom “looks”—does work, and so we used that model for the theory of the index of refraction.

18-3 The annihilation of positronium

We would like next to take an example which is very pretty. It is quite interesting and, although somewhat complicated, we hope not too much so. Our example is the system called *positronium*, which is an “atom” made up of an electron and a positron—a bound state of an e^+ and an e^- . It is like a hydrogen atom, except that a positron replaces the proton. This object has—like the hydrogen atom—many states. Also like the hydrogen, the ground state is split into a “hyperfine structure” by the interaction of the magnetic moments. The spins of the electron and positron are each one-half, and they can be either parallel or antiparallel to any given axis. (In the ground state there is no other angular momentum due to orbital motion.) So there are four states: three are the sub-states of a spin-one system, all with the same energy; and one is a state of spin zero with a different energy. The energy splitting is, however, much larger than the 1420 megacycles of hydrogen because the positron magnetic moment is so much stronger—1000 times stronger—than the proton moment.

The most important difference, however, is that positronium cannot last forever. The positron is the antiparticle of the electron; they can annihilate each other. The two particles disappear completely—converting their rest energy into radiation, which appears as γ -rays (photons). In the disintegration, two particles with a finite rest mass go into two or more objects which have zero rest mass.†

We begin by analyzing the disintegration of the spin-zero state of the positronium. It disintegrates into two γ -rays with a lifetime of about 10^{-10} second. Initially, we have a positron and an electron close together and with spins antiparallel, making the positronium system. After the disintegration there are two photons going out with equal and opposite momenta (Fig. 18-5). The momenta must be equal and opposite, because the total momentum after the disintegration must be zero, as it was before, if we are taking the case of annihilation at rest. If the positronium is not at rest, we can ride with it, solve the problem, and then transform everything back to the lab system. (See, we can do anything now; we have all the tools.)

First, we note that the angular distribution is not very interesting. Since the initial state has spin zero, it has no special axis—it is symmetric under all rotations. The final state must then also be symmetric under all rotations. That means that all angles for the disintegration are equally likely—the amplitude is the same for a photon to go in any direction. Of course, once we find *one* of the photons in some direction the *other* must be opposite.

† In the deeper understanding of the world today, we do not have an easy way to distinguish whether the energy of a photon is less “matter” than the energy of an electron, because as you remember all the particles behave very similarly. The only distinction is that the photon has zero rest mass.

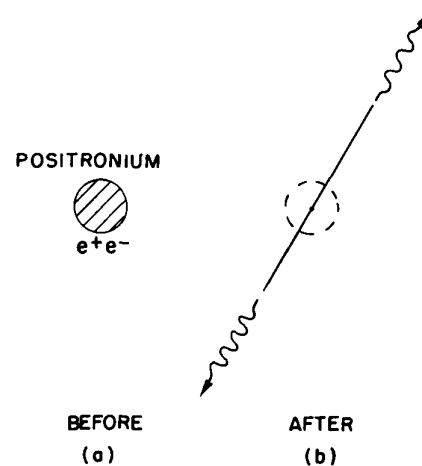


Fig. 18-5. The two-photon annihilation of positronium.

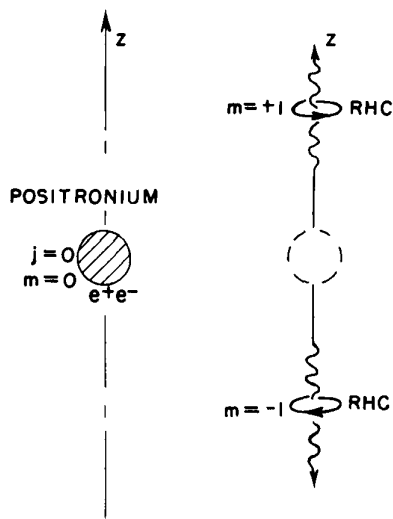


Fig. 18-6. One possibility for positronium annihilation along the z-axis.

The only remaining question, which we now want to look at, is about the polarization of the photons. Let's call the directions of motion of the two photons the plus and minus z-axes. We can use any representations we want for the polarization states of the photons; we will choose for our description right and left circular polarization—always with respect to the directions of motion. Right away, we can see that if the photon going upward is RHC, then angular momentum will be conserved if the downward going photon is also RHC. Each will carry $+1$ unit of angular momentum *with respect to its momentum direction*, which means plus and minus one unit about the z-axis. The total will be zero, and the angular momentum after the disintegration will be the same as before. See Fig. 18-6.

The same arguments show that if the upward going photon is RHC, the downward cannot be LHC. Then the final state would have two units of angular momentum. This is not permitted if the initial state has spin zero. Note that such a final state is also not possible for the other positronium ground state of spin one, because it can have a maximum of one unit of angular momentum in any direction.

Now we want to show that two-photon annihilation is not possible at all from the spin-one state. You might think that if we took the $j = 1, m = 0$ state—which has zero angular momentum about the z-axis—it should be like the spin-zero state, and could disintegrate into two RHC photons. Certainly, the disintegration sketched in Fig. 18-7(a) conserves angular momentum about the z-axis. But now look what happens if we rotate this system around the y-axis by 180° ; we get the picture shown in Fig. 18-7(b). It is exactly the same as in part (a) of the figure. All we have done is interchange the two photons. Now photons are Bose particles; if we interchange them, the amplitude has the same sign, so the amplitude for the disintegration in part (b) must be the same as in part (a). But we have assumed that the initial object is spin one. And when we rotate a spin-one object in a state with $m = 0$ by 180° about the y-axis, its amplitudes change sign (see Table 17-2 for $\theta = \pi$). So the amplitudes for (a) and (b) in Fig. 18-7 should have opposite signs; the spin-one state *cannot disintegrate into two photons*.

When positronium is formed you would expect it to end up in the spin-zero state $1/4$ of the time and in the spin-one state (with $m = -1, 0$, or $+1$) $3/4$ of the time. So $1/4$ of the time you would get two-photon annihilations. The other $3/4$

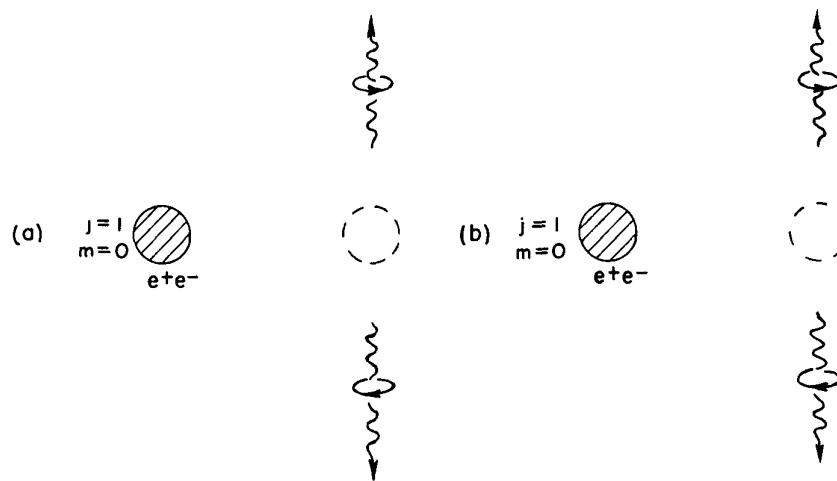


Fig. 18-7. For the $j = 1$ state of positronium, the process (a) and its 180° rotation about y (b) are exactly the same.

† Note that we always analyze the angular momentum about the direction of motion of the particle. If we were to ask about the angular momentum about any other axis, we would have to worry about the possibility of "orbital" angular momentum—from a $\mathbf{p} \times \mathbf{r}$ term. For instance, we can't say that the photons leave exactly from the center of the positronium. They could leave like two things shot out from the rim of a spinning wheel. We don't have to worry about such possibilities when we take our axis along the direction of motion.

of the time there can be no two-photon annihilations. There is still an annihilation, but it has to go with *three* photons. It is harder for it to do that and the lifetime is 1000 times longer—about 10^{-7} second. This is what is observed experimentally. We will not go into any more of the details of the spin-one annihilation.

So far we have that if we only worry about angular momentum, the spin-zero state of the positronium can go into two RHC photons. There is also another possibility: it can go into two LHC photons as shown in Fig. 18-8. The next question is, what is the relation between the amplitudes for these two possible decay modes? We can find out from the conservation of parity.

To do that, however, we need to know the parity of the positronium. Now theoretical physicists have shown in a way that is not easy to explain that the parity of the electron and the positron—its antiparticle—must be opposite, so that the spin-zero ground state of positronium must be odd. We will just assume that it is odd, and since we will get agreement with experiment, we can take that as sufficient proof.

Let's see then what happens if we make an inversion of the process in Fig. 18-6. When we do that, the two photons reverse directions *and* polarizations. The inverted picture looks just like Fig. 18-8. Assuming that the parity of the positronium is odd, the amplitudes for the two processes in Figs. 18-6 and 18-8 must have the opposite sign. Let's let $|R_1 R_2\rangle$ stand for the final state of Fig. 18-6 in which both photons are RHC, and let $|L_1 L_2\rangle$ stand for the final state of Fig. 18-8, in which both photons are LHC. The true final state—let's call it $|F\rangle$ —must be

$$|F\rangle = |R_1 R_2\rangle - |L_1 L_2\rangle. \quad (18.19)$$

Then an inversion changes the *R*'s into *L*'s and gives the state

$$P|F\rangle = |L_1 L_2\rangle - |R_1 R_2\rangle = -|F\rangle, \quad (18.20)$$

which is the negative of (18.19). So the final state $|F\rangle$ has negative parity, which is the same as the initial spin-zero state of the positronium. This is the only final state that conserves both angular momentum and parity. There is some amplitude that the disintegration into this state will occur, which we don't need to worry about now, however, since we are only interested in questions about the polarization.

What does the final state of (18.19) mean physically? One thing it means is the following: If we observe the two photons in two detectors which can be set to count separately the RHC or LHC photons, we will always see two RHC photons together, or two LHC photons together. That is, if you stand on one side of the positronium and someone else stands on the opposite side, you can measure the polarization and tell the other guy what polarization he will get. You have a 50-50 chance of catching a RHC photon or a LHC photon; whichever one you get, you can predict that he will get the same.

Since there is a 50-50 chance for RHC or LHC polarization, it sounds as though it might be like linear polarization. Let's ask what happens if we observe the photon in counters that accept only linearly polarized light. For γ -rays it is not as easy to measure the polarization as it is for light; there is no polarizer which works well for such short wavelengths. But let's imagine that there is, to make the discussion easier. Suppose that you have a counter that only accepts light with *x*-polarization, and that there is a guy on the other side that also looks for linear polarized light with, say, *y*-polarization. What is the chance you will pick up the two photons from an annihilation? What we need to ask is the amplitude that $|F\rangle$ will be in the state $|x_1 y_2\rangle$. In other words, we want the amplitude

$$\langle x_1 y_2 | F \rangle,$$

which is, of course, just

$$\langle x_1 y_2 | R_1 R_2 \rangle - \langle x_1 y_2 | L_1 L_2 \rangle. \quad (18.21)$$

Now although we are working with two-particle amplitudes for the two photons, we can handle them just as we did the single particle amplitudes, since

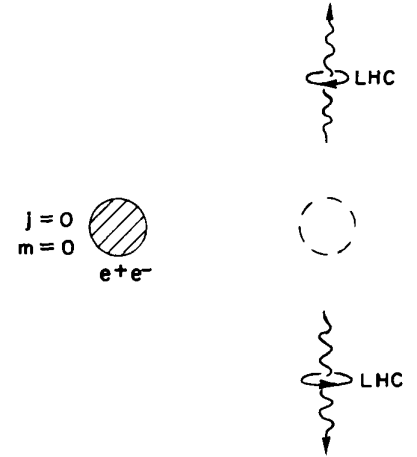


Fig. 18-8. Another possible process for positronium annihilation.

each particle acts independently of the other. That means that the amplitude $\langle x_1 y_2 | R_1 R_2 \rangle$ is just the product of the two independent amplitudes $\langle x_1 | R_1 \rangle$ and $\langle y_2 | R_2 \rangle$. Using Table 17-3, these two amplitudes are $1/\sqrt{2}$ and $i/\sqrt{2}$, so

$$\langle x_1 y_2 | R_1 R_2 \rangle = +\frac{i}{2}.$$

Similarly, we find that

$$\langle x_1 y_2 | L_1 L_2 \rangle = -\frac{i}{2}.$$

Subtracting these two amplitudes according to (18.21), we get that

$$\langle x_1 y_2 | F \rangle = +i. \quad (18.22)$$

So there is a *unit* probability† that if you get a photon in your *x*-polarized detector, the other guy will get a photon in his *y*-polarized detector.

Now suppose that the other guy sets his counter for *x*-polarization the same as yours. He would never get a count when you got one. If you work it through, you will find that

$$\langle x_1 x_2 | F \rangle = 0. \quad (18.23)$$

It will, naturally, also work out that if you set your counter for *y*-polarization he will get coincident counts only if he is set for *x*-polarization.

Now this all leads to an interesting situation. Suppose you were to set up something like a piece of calcite which separated the photons into *x*-polarized and *y*-polarized beams, and put a counter in each beam. Let's call one the *x*-counter and the other the *y*-counter. If the guy on the other side does the same thing, you can always tell him which beam his photon is going to go into. Whenever you and he get simultaneous counts, you can see which of your detectors caught the photon and then tell him which of his counters had a photon. Let's say that in a certain disintegration you find that a photon went into your *x*-counter; you can tell him that he must have had a count in his *y*-counter.

Now many people who learn quantum mechanics in the usual (old-fashioned) way find this disturbing. They would like to think that once the photons are emitted it goes along as a wave with a definite character. They would think that since "any given photon" has some "amplitude" to be *x*-polarized or to be *y*-polarized, there should be some chance of picking it up in either the *x*- or *y*-counter and that this chance shouldn't depend on what some other person finds out about a completely different photon. They argue that "someone else making a measurement shouldn't be able to change the probability that I will find something." Our quantum mechanics says, however, that by making a measurement on photon number one, you *can* predict precisely what the polarization of photon number two is going to be when it is detected. This point was never accepted by Einstein, and he worried about it a great deal—it became known as the "Einstein-Podolsky-Rosen paradox." But when the situation is described as we have done it here, there doesn't seem to be any paradox at all; it comes out quite naturally that what is measured in one place is correlated with what is measured somewhere else. The argument that the result is paradoxical runs something like this:

- (1) If you have a counter which tells you whether your photon is RHC or LHC, you can predict exactly what kind of a photon (RHC or LHC) he will find.
- (2) The photons he receives must, therefore, each be purely RHC or purely LHC, some of one kind and some of the other.
- (3) Surely you cannot alter the physical nature of *his* photons by changing the kind of observation you make on *your* photons. No matter what measurements you make on yours, his must still be either RHC or LHC.

† We have not normalized our amplitudes, or multiplied them by the amplitude for the disintegration into any particular final state, but we can see that this result is correct because we get zero probability when we look at the other alternative—see Eq. (18.23).

- (4) Now suppose he changes his apparatus to split his photons into two linearly polarized beams with a piece of calcite so that all of his photons go either into an x -polarized beam or into a y -polarized beam. There is absolutely no way, according to quantum mechanics, to tell into which beam any particular RHC photon will go. There is a 50% probability it will go into the x -beam and a 50% probability it will go into the y -beam. And the same goes for a LHC photon.
- (5) Since each photon is RHC or LHC—according to (2) and (3)—each one must have a 50–50 chance of going into the x -beam or the y -beam and there is no way to predict which way it will go.
- (6) Yet the theory predicts that if *you* see your photon go through an x -polarizer you can predict *with certainty* that his photon will go into his y -polarized beam. This is in contradiction to (5) so there is a paradox.

Nature apparently doesn't see the "paradox," however, because experiment shows that the prediction in (6) is, in fact, true. We have already discussed the key to this "paradox" in our very first lecture on quantum mechanical behavior in Chapter 35, Vol. I. In the argument above, steps (1), (2), (4), and (6) are all correct, but (3), and its consequence (5), are wrong; they are not a true description of nature. Argument (3) says that by *your* measurement (seeing a RHC or a LHC photon) you can determine which of two alternative events occurs for him (seeing a RHC or a LHC photon), *and* that even if you do *not* make your measurement you can still say that his event will occur either by one alternative or the other. But it was precisely the point of Chapter 35, Vol. I, to point out right at the beginning that this is not so in Nature. *Her* way requires a description in terms of interfering amplitudes, one amplitude for each alternative. A measurement of which alternative actually occurs destroys the interference, but if a measurement is *not* made you cannot still say that "one alternative or the other is still occurring."

If you could determine for each one of your photons whether it was RHC and LHC, and *also* whether it was x -polarized (all for the same photon) there would indeed be a paradox. But you cannot do that—it is an example of the uncertainty principle.

Do you still think there is a "paradox"? Make sure that it is, in fact, a paradox about the behavior of Nature, by setting up an imaginary experiment for which the theory of quantum mechanics would predict inconsistent results via two different arguments. Otherwise the "paradox" is only a conflict between reality and your feeling of what reality "ought to be."

Do you think that it is *not* a "paradox," but that it is still very peculiar? On that we can all agree. It is what makes physics fascinating.

18-4 Rotation matrix for any spin

By now you can see, we hope, how important the idea of the angular momentum is in understanding atomic processes. So far, we have considered only systems with spins—or "total angular momentum"—of zero, one-half, or one. There are, of course, atomic systems with higher angular momenta. For analyzing such systems we would need to have tables of rotation amplitudes like those in Section 17-6. That is, we would need the matrix of amplitudes for spin $\frac{3}{2}$, 2, $\frac{5}{2}$, 3, etc. Although we will not work out these tables in detail, we would like to show you how it is done, so that you can do it if you ever need to.

As we have seen earlier, any system which has the spin or "total angular momentum" j can exist in any one of $(2j + 1)$ states for which the z -component of angular momentum can have any one of the discrete values in the sequence $j, j - 1, j - 2, \dots, -(j - 1), -j$ (all in units of \hbar). Calling the z -component of angular momentum of any particular state $m\hbar$, we can define a particular angular momentum state by giving the numerical values of the two "angular momentum quantum numbers" j and m . We can indicate such a state by the state vector $|j, m\rangle$. In the case of a spin one-half particle, the two states are then $|\frac{1}{2}, \frac{1}{2}\rangle$ and $|\frac{1}{2}, -\frac{1}{2}\rangle$; or for a spin-one system, the states would be written in this notation as $|1, +1\rangle, |1, 0\rangle, |1, -1\rangle$. A spin-zero particle has, of course, only the one state $|0, 0\rangle$.

Now we want to know what happens when we project the general state $|j, m\rangle$ into a representation with respect to a rotated set of axes. First, we know that j is a number which characterizes *the system*, so it doesn't change. If we rotate the axes, all we do is get a mixture of the various m -values for the same j . In general, there will be some amplitude that in the rotated frame the system will be in the state $|j, m'\rangle$, where m' gives the new z -component of angular momentum. So what we want are all the matrix elements $\langle j, m' | R | j, m \rangle$ for various rotations. We already know what happens if we rotate by an angle ϕ about the z -axis. The new state is just the old one multiplied by $e^{im\phi}$ —it still has the same m -value. We can write this by

$$R_z(\phi) |j, m\rangle = e^{im\phi} |j, m\rangle. \quad (18.24)$$

Or, if you prefer,

$$\langle j, m' | R_z(\phi) | j, m \rangle = \delta_{m, m'} e^{im\phi} \quad (18.25)$$

(where $\delta_{m, m'}$ is 1 if $m' = m$, or zero otherwise).

For a rotation about any other axis there will be a mixing of the various m -states. We could, of course, try to work out the matrix elements for an arbitrary rotation described by the Euler angles β , α , and γ . But it is easier to remember that the most general such rotation can be made up of the three rotations $R_z(\gamma)$, $R_y(\alpha)$, $R_z(\beta)$; so if we know the matrix elements for a rotation about the y -axis, we will have all we need.

How can we find the rotation matrix for a rotation by the angle θ about the y -axis for a particle of spin j ? We can't tell you how to do it in a basic way (with what we have had). We did it for spin one-half by a complicated symmetry argument. We then did it for spin one by taking the special case of a spin-one system which was made up of two spin one-half particles. If you will go along with us and accept the fact that in the general case the answers depend only on the spin j , and are independent of how the inner guts of the object of spin j are put together, we can extend the spin-one argument to an arbitrary spin. We can, for example, cook up an artificial system of spin $\frac{3}{2}$ out of three spin one-half objects. We can even avoid complications by imagining that they are all distinct particles—like a proton, an electron, and a muon. By transforming each spin one-half object, we can see what happens to the whole system—remembering that the three amplitudes are multiplied for the combined state. Let's see how it goes in this case.

Suppose we take the three spin one-half objects all with spins “up”; we can indicate this state by $|+++ \rangle$. If we look at this system in a frame rotated about the z -axis by the angle ϕ , each plus stays a plus, but gets multiplied by $e^{i\phi/2}$. We have three such factors, so

$$R_z(\phi) |+++ \rangle = e^{i(3\phi/2)} |+++ \rangle. \quad (18.26)$$

Evidently the state $|+++ \rangle$ is just what we mean by the $m = +\frac{3}{2}$ state, or the state $|\frac{3}{2}, +\frac{3}{2}\rangle$.

If we now rotate this system about the y -axis, each of the spin one-half objects will have some amplitude to be plus or to be minus, so the system will now be a mixture of the *eight* possible combinations $|+++ \rangle$, $|++- \rangle$, $|+-+ \rangle$, $|-++ \rangle$, $|+- - \rangle$, $|-+ - \rangle$, $|- - + \rangle$, or $|- - - \rangle$. It is clear, however, that these can be broken up into four sets, each set corresponding to a particular value of m . First, we have $|+++ \rangle$, for which $m = \frac{3}{2}$. Then there are the three states $|++- \rangle$, $|+-+ \rangle$, and $|-++ \rangle$ —each with two pluses and one minus. Since each spin one-half object has the same chance of coming out minus under the rotation, the amounts of each of these three combinations should be equal. So let's take the combination

$$\frac{1}{\sqrt{3}} \{ |++- \rangle + |+-+ \rangle + |-++ \rangle \} \quad (18.27)$$

with the factor $1/\sqrt{3}$ put in to normalize the state. If we rotate this state about the z -axis, we get a factor $e^{i\phi/2}$ for each plus, and $e^{-i\phi/2}$ for each minus. Each term in (18.27) is multiplied by $e^{i\phi/2}$, so there is the common factor $e^{i\phi/2}$. This

one “—” pieces. For instance,

$$\begin{aligned} |++-\rangle &= a^2c|+'+'+\rangle + a^2d|+'+'-\rangle + abc|+'-+'+\rangle \\ &\quad + bac|+'-+'+\rangle + abd|+'-+'-\rangle + bad|+'-+'-\rangle \\ &\quad + b^2c|+'-+'+\rangle + b^2d|+'-+'-\rangle. \end{aligned} \quad (18.33)$$

Adding two similar expressions for $|+-+\rangle$ and $| -++\rangle$ and dividing by $\sqrt{3}$, we find

$$\begin{aligned} |\frac{3}{2}, +\frac{1}{2}, S\rangle &= \sqrt{3} a^2c |\frac{3}{2}, +\frac{3}{2}, T\rangle \\ &\quad + (a^2d + 2abc) |\frac{3}{2}, +\frac{1}{2}, T\rangle \\ &\quad + (2bad + b^2c) |\frac{3}{2}, -\frac{1}{2}, T\rangle \\ &\quad + \sqrt{3} b^2d |\frac{3}{2}, -\frac{3}{2}, T\rangle. \end{aligned} \quad (18.34)$$

Continuing the process we find all the elements $\langle jT | iS \rangle$ of the transformation matrix as given in Table 18-2. The first column comes from Eq. (18.32); the second from (18.34). The last two columns were worked out in the same way.

Table 18-2

Rotation matrix for a spin $\frac{3}{2}$ particle

(The coefficients a , b , c , and d are given in Table 12-4.)

$\langle jT iS \rangle$	$ \frac{3}{2}, +\frac{3}{2}, S\rangle$	$ \frac{3}{2}, +\frac{1}{2}, S\rangle$	$ \frac{3}{2}, -\frac{1}{2}, S\rangle$	$ \frac{3}{2}, -\frac{3}{2}, S\rangle$
$\langle \frac{3}{2}, +\frac{3}{2}, T $	a^3	$\sqrt{3} a^2c$	$\sqrt{3} ac^2$	c^3
$\langle \frac{3}{2}, +\frac{1}{2}, T $	$\sqrt{3} a^2b$	$a^2d + 2abc$	$c^2b + 2dac$	$\sqrt{3} c^2d$
$\langle \frac{3}{2}, -\frac{1}{2}, T $	$\sqrt{3} ab^2$	$2bad + b^2c$	$2cdb + d^2a$	$\sqrt{3} cd^2$
$\langle \frac{3}{2}, -\frac{3}{2}, T $	b^3	$\sqrt{3} b^2d$	$\sqrt{3} bd^2$	d^3

Now suppose the T -frame were rotated with respect to S by the angle θ about their y -axes. Then a , b , c , and d have the values [see (12.54)] $a = d = \cos \theta/2$, and $c = -b = \sin \theta/2$. Using these values in Table 18-2 we get the forms which correspond to the second part of Table 17-2, but now for a spin $\frac{3}{2}$ system.

The arguments we have just gone through are readily generalized to a system of any spin j . The states $|j, m\rangle$ can be put together from $2j$ particles, each of spin one-half. (There are $j + m$ of them in the $|+\rangle$ state and $j - m$ in the $|-\rangle$ state.) Sums are taken over all the possible ways this can be done, and the state is normalized by multiplying by a suitable constant. Those of you who are mathematically inclined may be able to show that the following result comes out†:

$$\begin{aligned} \langle j, m' | R_y(\theta) | j, m \rangle &= [(j + m)!(j - m)!(j + m')!(j - m')!]^{1/2} \\ &\quad \times \sum_k \frac{(-1)^k (\cos \theta/2)^{2j+m'-m-2k} (\sin \theta/2)^{m-m'+2k}}{(m - m' + k)!(j + m' - k)!(j - m - k)!k!}, \end{aligned} \quad (18.35)$$

where k is to go over all values which give terms ≥ 0 in all the factorials.

This is quite a messy formula, but with it you can check Table 17-2 for $j = 1$ and prepare tables of your own for larger j . Several special matrix elements are of extra importance and have been given special names. For example the matrix elements for $m = m' = 0$ and integral j are known as the Legendre polynomials and are called $P_j(\cos \theta)$:

$$\langle j, 0 | R_y(\theta) | j, 0 \rangle = P_j(\cos \theta). \quad (18.36)$$

† If you want details, they are given in an appendix to this chapter.

The first few of these polynomials are:

$$P_0(\cos \theta) = 1, \quad (18.37)$$

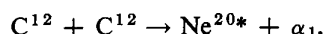
$$P_1(\cos \theta) = \cos \theta, \quad (18.38)$$

$$P_2(\cos \theta) = \frac{1}{2}(3 \cos^2 \theta - 1), \quad (18.39)$$

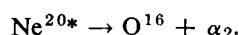
$$P_3(\cos \theta) = \frac{1}{2}(5 \cos^3 \theta - 3 \cos \theta). \quad (18.40)$$

18-5 Measuring a nuclear spin

We would like to show you one example of the application of the coefficients we have just described. It has to do with a recent, interesting experiment which you will now be able to understand. Some physicists wanted to find out the spin of a certain excited state of the Ne^{20} nucleus. To do this, they bombarded a carbon target with a beam of accelerated carbon ions, and produced the desired excited state of Ne^{20} —called Ne^{20*} —in the reaction



where α_1 is the α -particle, or He^4 . Several of the excited states of Ne^{20} produced this way are unstable and disintegrate in the reaction



So experimentally there are two α -particles which come out of the reaction. We call them α_1 and α_2 ; since they come off with different energies, they can be distinguished from each other. Also, by picking a particular energy for α_1 we can pick out any particular excited state of the Ne^{20} .

The experiment was set up as shown in Fig. 18-9. A beam of 16-Mev carbon ions was directed onto a thin foil of carbon. The first α -particle was counted in a silicon diffused junction detector marked α_1 —set to accept α -particles of the proper energy moving in the forward direction (with respect to the incident C^{12} beam). The second α -particle was picked up in the counter α_2 at the angle θ with respect to α_1 . The counting rate of coincidence signals from α_1 and α_2 were measured as a function of the angle θ .

The idea of the experiment is the following. First, you need to know that the spins of C^{12} , O^{16} , and the α -particle are all zero. If we call the direction of motion of the initial C^{12} the $+z$ -direction, then we know that the Ne^{20*} must have zero angular momentum about the z -axis. None of the other particles has any spin; the C^{12} arrives along the z -axis and the α_1 leaves along the z -axis so they can't have any angular momentum about it. So whatever the spin j of the Ne^{20*} is, we know that it is in the state $|j, 0\rangle$. Now what will happen when the Ne^{20*} disintegrates into an O^{16} and the second α -particle? Well, the α -particle is picked up in the counter α_2 and to conserve momentum the O^{16} must go off in the opposite direction.† *About the new axis* through α_2 , there can be no component of angular momentum. The final state has zero angular momentum about the new axis, so the Ne^{20*} can disintegrate this way only if it has some amplitude to have m' equal to zero, where m' is the quantum number of the component of angular momentum about the new axis. In fact, the probability of observing α_2 at the angle θ is just the square of the amplitude (or matrix element)

$$\langle j, 0 | R_y(\theta) | j, 0 \rangle. \quad (18.41)$$

To find the spin of the Ne^{20*} state in question, the intensity of the second α -particle was plotted as a function of angle and compared with the theoretical

† We can neglect the recoil given to the Ne^{20*} in the first collision. Or better still, we can calculate what it is and make a correction for it.

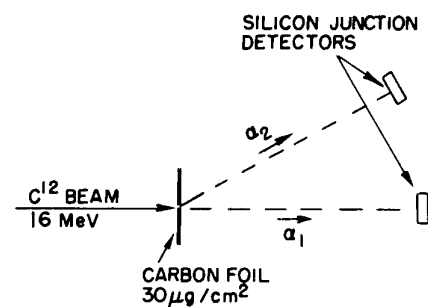


Fig. 18-9. Experimental arrangement for determining the spin of certain states of Ne^{20} .

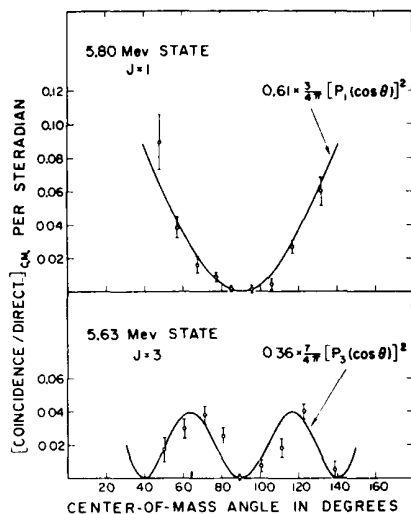


Fig. 18-10. Experimental results for the angular distribution of the α -particles from two excited states of Ne^{20} produced in the setup of Fig. 18-9. [From J. A. Kuehner, *Physical Review*, Vol. 125, p. 1653, 1962.]

curves for various values of j . As we said in the last section, the amplitudes $\langle j, 0 | R_y(\theta) | j, 0 \rangle$ are just the functions $P_j(\cos \theta)$. So the possible angular distributions are curves of $[P_j(\cos \theta)]^2$. The experimental results are shown in Fig. 18-10 for two of the excited states. You can see that the angular distribution for the 5.80-Mev state fits very well the curve for $[P_1(\cos \theta)]^2$, and so it must be a spin-one state. The data for the 5.63-Mev state, on the other hand, are quite different; they fit the curve $[P_3(\cos \theta)]^2$. The state has a spin of 3.

From this experiment we have been able to find out the angular momentum of two of the excited states of Ne^{20*} . This information can then be used for trying to understand what the configuration of protons and neutrons is inside this nucleus—one more piece of information about the mysterious nuclear forces.

18-6 Composition of angular momentum

When we studied the hyperfine structure of the hydrogen atom in Chapter 12 we had to work out the internal states of a system composed of two particles—the electron and the proton—each with a spin of one-half. We found that the four possible spin states of such a system could be put together into two groups—a group with one energy that looked to the external world like a spin-one particle, and one remaining state that behaved like a particle of zero spin. That is, putting together two spin one-half particles we can form a system whose “total spin” is one, or zero. In this section we want to discuss in more general terms the spin states of a *system* which is made up of two particles of arbitrary spin. It is another important problem about angular momentum in quantum mechanical systems.

Let's first rewrite the results of Chapter 12 for the hydrogen atom in a form that will be easier to extend to the more general case. We began with two particles which we will now call particle a (the electron) and particle b (the proton). Particle a had the spin $j_a (= \frac{1}{2})$, and its z -component of angular momentum m_a could have one of several values (actually 2, namely $m_a = +\frac{1}{2}$ or $m_a = -\frac{1}{2}$). Similarly, the spin state of particle b is described by its spin j_b and its z -component of angular momentum m_b . Various combinations of the spin states of the two particles could be formed. For instance, we could have particle a with $m_a = \frac{1}{2}$ and particle b with $m_b = -\frac{1}{2}$, to make a state $|a, +\frac{1}{2}; b, -\frac{1}{2}\rangle$. In general, the combined states formed a system whose “system spin,” or “total spin,” or “total angular momentum” J could be 1, or 0. And the system could have a z -component of angular momentum M , which was $+1$, 0, or -1 when $J = 1$, or 0 when $J = 0$. In this new language we can rewrite the formulas in (12.41) and (12.42) as shown in Table 18-3.

In the table the left-hand column describes the compound state in terms of its total angular momentum J and the z -component M . The right-hand column shows how these states are made up in terms of the m -values of the two particles a and b .

We want now to generalize this result to states made up of two objects a and b of arbitrary spins j_a and j_b . We start by considering an example for which $j_a = \frac{1}{2}$

Table 18-3

Composition of angular momenta for two spin $\frac{1}{2}$ particles ($j_a = \frac{1}{2}, j_b = \frac{1}{2}$)

$ J = 1, M = +1\rangle$	$= a, +\frac{1}{2}; b, +\frac{1}{2}\rangle$
$ J = 1, M = 0\rangle$	$= \frac{1}{\sqrt{2}} \{ a, +\frac{1}{2}; b, -\frac{1}{2}\rangle + a, -\frac{1}{2}; b, +\frac{1}{2}\rangle \}$
$ J = 1, M = -1\rangle$	$= a, -\frac{1}{2}; b, -\frac{1}{2}\rangle$
$ J = 0, M = 0\rangle$	$= \frac{1}{\sqrt{2}} \{ a, +\frac{1}{2}; b, -\frac{1}{2}\rangle - a, -\frac{1}{2}; b, +\frac{1}{2}\rangle \}$

and $j_b = 1$, namely, the deuterium atom in which particle a is an electron (e) and particle b is the nucleus—a deuteron (d). We have then that $j_a = j_e = \frac{1}{2}$. The deuteron is formed of one proton and one neutron in a state whose total spin is one, so $j_b = j_d = 1$. We want to discuss the hyperfine states of deuterium—just the way we did for hydrogen. Since the deuteron has three possible states $m_b = m_d = +1, 0, -1$, and the electron has two, $m_a = m_e = +\frac{1}{2}, -\frac{1}{2}$, there are six possible states as follows (using the notation $|e, m_e; d, m_d\rangle$):

$$\begin{aligned} & |e, +\frac{1}{2}; d, +1\rangle, \\ & |e, +\frac{1}{2}; d, 0\rangle; |e, -\frac{1}{2}; d, +1\rangle, \\ & |e, +\frac{1}{2}; d, -1\rangle; |e, -\frac{1}{2}; d, 0\rangle, \\ & |e, -\frac{1}{2}; d, -1\rangle. \end{aligned} \quad (18.42)$$

You will notice that we have grouped the states according to the values of the sum of m_e and m_d —arranged in descending order.

Now we ask: What happens to these states if we project into a different coordinate system? If the new system is just rotated about the z -axis by the angle ϕ , then the state $|e, m_e; d, m_d\rangle$ gets multiplied by

$$e^{im_e\phi} e^{im_d\phi} = e^{i(m_e+m_d)\phi}. \quad (18.43)$$

(The state may be thought of as the product $|e, m_e\rangle |d, m_d\rangle$, and each state vector contributes independently its own exponential factor.) The factor (18.43) is of the form $e^{iM\phi}$, so the state $|e, m_e; d, m_d\rangle$ has a z -component of angular momentum equal to

$$M = m_e + m_d. \quad (18.44)$$

The z -component of the total angular momentum is the sum of the z -components of angular momentum of the parts.

In the list of (18.42), therefore, the state in the top line has $M = +\frac{3}{2}$, the two in the second line have $M = +\frac{1}{2}$, the next two have $M = -\frac{1}{2}$, and the last state has $M = -\frac{3}{2}$. We see immediately one possibility for the spin J of the combined state (the total angular momentum) must be $\frac{3}{2}$, and this will require four states with $M = +\frac{3}{2}, +\frac{1}{2}, -\frac{1}{2}$, and $-\frac{3}{2}$.

There is only one candidate for $M = \frac{3}{2}$, so we know already that

$$|J = \frac{3}{2}, M = +\frac{3}{2}\rangle = |e, +\frac{1}{2}; d, +1\rangle. \quad (18.45)$$

But what is the state $|J = \frac{3}{2}, M = \frac{1}{2}\rangle$? We have two candidates in the second line of (18.42), and, in fact, any linear combination of them would also have $M = \frac{1}{2}$. So, in general, we must expect to find that

$$|J = \frac{3}{2}, M = +\frac{1}{2}\rangle = \alpha |e, +\frac{1}{2}; d, 0\rangle + \beta |e, -\frac{1}{2}; d, +1\rangle, \quad (18.46)$$

where α and β are two numbers. They are called the *Clebsch-Gordon coefficients*. Our next problem is to find out what they are.

We can find out easily if we just remember that the deuteron is made up of a neutron and a proton, and write the deuteron states out more explicitly using the rules of Table 18-3. If we do that, the states listed in (18.42) then look as shown in Table 18-4.

We want to form the four states of $J = \frac{3}{2}$, using the states in the table. But we already know the answer, because in Table 18-1 we have states of spin $\frac{3}{2}$ formed from three spin one-half particles. The first state in Table 18-1 has $|J = \frac{3}{2}, M = +\frac{3}{2}\rangle$ and it is $|+++\rangle$, which—in our present notation—is the same as $|e, +\frac{1}{2}; n, +\frac{1}{2}, p, +\frac{1}{2}\rangle$, or the first state in Table 18-4. But this state is also the same as the first in the list of (18.42), confirming our statement in (18.45). The second line of Table 18-1 says—changing to our present notation—that

$$\begin{aligned} |J = \frac{3}{2}, M = +\frac{1}{2}\rangle &= \frac{1}{\sqrt{3}} \{ |e, +\frac{1}{2}; n, +\frac{1}{2}; p, -\frac{1}{2}\rangle \\ &+ |e, +\frac{1}{2}; n, -\frac{1}{2}; p, +\frac{1}{2}\rangle + |e, -\frac{1}{2}; n, +\frac{1}{2}; p, +\frac{1}{2}\rangle \}. \end{aligned} \quad (18.47)$$

Table 18-4

Angular momentum states for a deuterium atom

$m = \frac{3}{2}$ $ e, +\frac{1}{2}; d, +1\rangle = e, +\frac{1}{2}; n, +\frac{1}{2}; p, +\frac{1}{2}\rangle$
$m = \frac{1}{2}$ $ e, +\frac{1}{2}; d, 0\rangle = \frac{1}{\sqrt{2}} \{ e, +\frac{1}{2}; n, +\frac{1}{2}; p, -\frac{1}{2}\rangle + e, +\frac{1}{2}; n, -\frac{1}{2}; p, +\frac{1}{2}\rangle \}$ $ e, -\frac{1}{2}; d, +1\rangle = e, -\frac{1}{2}; n, +\frac{1}{2}; p, +\frac{1}{2}\rangle$
$m = -\frac{1}{2}$ $ e, +\frac{1}{2}; d, -1\rangle = e, +\frac{1}{2}; n, -\frac{1}{2}; p, -\frac{1}{2}\rangle$ $ e, -\frac{1}{2}; d, 0\rangle = \frac{1}{\sqrt{2}} \{ e, -\frac{1}{2}; n, +\frac{1}{2}; p, -\frac{1}{2}\rangle + e, -\frac{1}{2}; n, -\frac{1}{2}; p, +\frac{1}{2}\rangle \}$
$m = -\frac{3}{2}$ $ e, -\frac{1}{2}; d, -1\rangle = e, -\frac{1}{2}; n, -\frac{1}{2}; p, -\frac{1}{2}\rangle$

The right side can evidently be put together from the two entries in the second line of Table 18-4 by taking $\sqrt{2/3}$ of the first term with $\sqrt{1/3}$ of the second. That is, Eq. (18-47) is equivalent to

$$|J = \frac{3}{2}, M = \frac{1}{2}\rangle = \sqrt{2/3} |e, +\frac{1}{2}; d, 0\rangle + \sqrt{1/3} |e, -\frac{1}{2}; d, 1\rangle. \quad (18.48)$$

We have found our two Clebsch-Gordon coefficients α and β in Eq. (18.46):

$$\alpha = \sqrt{2/3}, \quad \beta = \sqrt{1/3}. \quad (18.49)$$

Following the same procedure we can find that

$$|J = \frac{3}{2}, M = -\frac{1}{2}\rangle = \sqrt{1/3} |e, +\frac{1}{2}; d, -1\rangle + \sqrt{2/3} |e, -\frac{1}{2}; d, 0\rangle. \quad (18.50)$$

And, also, of course,

$$|J = \frac{3}{2}, M = -\frac{3}{2}\rangle = |e, -\frac{1}{2}; d, -1\rangle. \quad (18.51)$$

These are the rules for the composition of spin 1 and spin $\frac{1}{2}$ to make a total $J = \frac{3}{2}$. We summarize (18.45), (18.48), and (18.50) in Table 18-5.

We have, however, only four states here while the system we are considering has six possible states. Of the two states in the second line of (18.42) we have used only one linear combination to form $|J = \frac{3}{2}, M = +\frac{1}{2}\rangle$. There is another linear combination orthogonal to the one we have taken which also has $M = +\frac{1}{2}$, namely

$$\sqrt{1/3} |e, +\frac{1}{2}; d, 0\rangle - \sqrt{2/3} |e, -\frac{1}{2}; d, +1\rangle. \quad (18.52)$$

Table 18-5

The $J = \frac{3}{2}$ states of the deuterium atom

$ J = \frac{3}{2}, M = +\frac{3}{2}\rangle = e, +\frac{1}{2}; d, +1\rangle$ $ J = \frac{3}{2}, M = +\frac{1}{2}\rangle = \sqrt{2/3} e, +\frac{1}{2}; d, 0\rangle + \sqrt{1/3} e, -\frac{1}{2}; d, 1\rangle$ $ J = \frac{3}{2}, M = -\frac{1}{2}\rangle = \sqrt{1/3} e, +\frac{1}{2}; d, -1\rangle + \sqrt{2/3} e, -\frac{1}{2}; d, 0\rangle$ $ J = \frac{3}{2}, M = -\frac{3}{2}\rangle = e, -\frac{1}{2}; d, -1\rangle$
--

Similarly, the two states in the third line of (18.42) can be combined to give two orthogonal states, each with $M = -\frac{1}{2}$. The one orthogonal to (18.52) is

$$\sqrt{2/3} |e, +\frac{1}{2}; d, -1\rangle - \sqrt{1/3} |e, -\frac{1}{2}; d, 0\rangle. \quad (18.53)$$

These are the two remaining states. They have $M = m_e + m_d = \pm\frac{1}{2}$; and must be the two states corresponding to $J = \frac{1}{2}$. So we have

$$\begin{aligned} |J = \frac{1}{2}, M = \frac{1}{2}\rangle &= \sqrt{1/3} |e, +\frac{1}{2}; d, 0\rangle - \sqrt{2/3} |e, -\frac{1}{2}; d, +1\rangle, \\ |J = \frac{1}{2}, M = -\frac{1}{2}\rangle &= \sqrt{2/3} |e, +\frac{1}{2}; d, -1\rangle - \sqrt{1/3} |e, -\frac{1}{2}; d, 0\rangle. \end{aligned} \quad (18.54)$$

We can verify that these two states do indeed behave like the states of a spin one-half object by writing out the deuteron parts in terms of the neutron and proton states—using Table 18-3. The first state in (18.53) is

$$\begin{aligned} \sqrt{1/6} \{ &|e, +\frac{1}{2}; n, +\frac{1}{2}; p, -\frac{1}{2}\rangle + |e, +\frac{1}{2}; n, -\frac{1}{2}; p, +\frac{1}{2}\rangle \\ &- \sqrt{2/3} |e, -\frac{1}{2}; n, +\frac{1}{2}; p, +\frac{1}{2}\rangle, \end{aligned} \quad (18.55)$$

which can also be written

$$\begin{aligned} \sqrt{1/3} [&\sqrt{1/2} \{ |e, +\frac{1}{2}; n, +\frac{1}{2}; p, -\frac{1}{2}\rangle - |e, -\frac{1}{2}; n, +\frac{1}{2}; p, +\frac{1}{2}\rangle \\ &+ \sqrt{1/2} \{ |e, +\frac{1}{2}; n, -\frac{1}{2}; p, +\frac{1}{2}\rangle - |e, -\frac{1}{2}; n, +\frac{1}{2}; p, +\frac{1}{2}\rangle \}]. \end{aligned} \quad (18.56)$$

Now look at the terms in the first curly brackets, and think of the e and p taken together. Together they form a spin-zero state (see the bottom line of Table 18-3), and contribute no angular momentum. Only the neutron is left, so the whole of the *first* curly bracket of (18.56) behaves under rotations like a neutron, namely as a state with $J = \frac{1}{2}$, $M = +\frac{1}{2}$. Following the same reasoning, we see that in the *second* curly bracket of (18.56) the electron and neutron team up to produce zero angular momentum, and only the proton contribution—with $m_p = \frac{1}{2}$ —is left. The terms behave like an object with $J = \frac{1}{2}$, $M = +\frac{1}{2}$. So the whole expression of (18.56) transforms like $|J = +\frac{1}{2}, M = +\frac{1}{2}\rangle$ as it should. The $M = -\frac{1}{2}$ state which corresponds to (18.57) can be written down (by changing the proper $+\frac{1}{2}$'s to $-\frac{1}{2}$'s) to get

$$\begin{aligned} \sqrt{1/3} [&\sqrt{1/2} \{ |e, +\frac{1}{2}; n, -\frac{1}{2}; p, -\frac{1}{2}\rangle - |e, -\frac{1}{2}; n, -\frac{1}{2}; p, +\frac{1}{2}\rangle \\ &+ \sqrt{1/2} \{ |e, +\frac{1}{2}; n, -\frac{1}{2}; p, -\frac{1}{2}\rangle - |e, -\frac{1}{2}; n, +\frac{1}{2}; p, -\frac{1}{2}\rangle \}]. \end{aligned} \quad (18.57)$$

You can easily check that this is equal to the second line of (18.54), as it should be if the two terms of that pair are to be the two states of a spin one-half system. So our results are confirmed. A deuteron and an electron can exist in six spin states, four of which act like the states of a spin $\frac{3}{2}$ object (Table 18-5) and two of which act like an object of spin one-half (18.54).

The results of Table 18-5 and of Eq. (18.54) were obtained by making use of the fact that the deuteron is made up of a neutron and a proton. The truth of the equations does not depend on that special circumstance. For *any* spin-one object put together with any spin one-half object the composition laws (and the coefficients) are the same. The set of equations in Table 18-5 means that if the coordinates are rotated about, say, the y-axis—so that the states of the spin one-half particle and of the spin-one particle change according to Table 18-1 and Table 18-2—the linear combinations on the right-hand side will change in the proper way for a spin $\frac{3}{2}$ object. Under the same rotation the states of (18.54) will change as the states of a spin one-half object. The results depend only on the

Table 18-6

Composition of a spin one-half particle ($j_a = \frac{1}{2}$)
and a spin-one particle ($j_b = 1$).

$ J = \frac{3}{2}, M = \frac{3}{2}\rangle = a, +\frac{1}{2}; b, +1\rangle$
$ J = \frac{3}{2}, M = \frac{1}{2}\rangle = \sqrt{2/3} a, +\frac{1}{2}; b, 0\rangle + \sqrt{1/3} a, -\frac{1}{2}; b, +1\rangle$
$ J = \frac{3}{2}, M = -\frac{1}{2}\rangle = \sqrt{1/3} a, +\frac{1}{2}; b, -1\rangle + \sqrt{2/3} a, -\frac{1}{2}; b, 0\rangle$
$ J = \frac{3}{2}, M = -\frac{3}{2}\rangle = a, -\frac{1}{2}; b, -1\rangle$
$ J = \frac{1}{2}, M = +\frac{1}{2}\rangle = \sqrt{1/3} a, +\frac{1}{2}; b, 0\rangle - \sqrt{2/3} a, -\frac{1}{2}; b, +1\rangle$
$ J = \frac{1}{2}, M = -\frac{1}{2}\rangle = \sqrt{2/3} a, +\frac{1}{2}; b, -1\rangle - \sqrt{1/3} a, -\frac{1}{2}; b, 0\rangle$

rotation properties (that is, the spin states) of the two original particles but not in any way on the origins of their angular momenta. We have only made use of this fact to work out the formulas by choosing a special case in which one of the component parts is itself made up of two spin one-half particles in a symmetric state. We have put all our results together in Table 18-6, changing the notation "c" and "d" to "a" and "b" to emphasize the generality of the conclusions.

Suppose we have the general problem of finding the states which can be formed when two objects of arbitrary spins are combined. Say one has j_a (so its z-component m_a runs over the $2j_a + 1$ values from $-j_a$ to $+j_a$) and the other has j_b (with z-component m_b running over the values from $-j_b$ to $+j_b$). The combined states are $|a, m_a; b, m_b\rangle$, and there are $(2j_a + 1)(2j_b + 1)$ different ones. Now what states of total spin J can be found?

The total z-component of angular momentum M is equal to $m_a + m_b$, and the states can all be listed according to M [as in (18.42)]. The largest M is unique; it corresponds to $m_a = j_a$ and $m_b = j_b$, and is, therefore, just $j_a + j_b$. That means that the largest total spin J is also equal to the sum $j_a + j_b$:

$$J = (M)_{\max} = j_a + j_b.$$

For the first M value smaller than $(M)_{\max}$, there are two states (either m_a or m_b is one unit less than its maximum). They must contribute one state to the set that goes with $J = j_a + j_b$, and the one left over will belong to a new set with $J = j_a + j_b - 1$. The next M -value—the third from the top of the list—can be formed in *three* ways. (From $m_a = j_a - 2, m_b = j_b$; from $m_a = j_a - 1, m_b = j_b - 1$; and from $m_a = j_a, m_b = j_b - 2$.) Two of these belong to groups already started above; the third tells us that states of $J = j_a + j_b - 2$ must also be included. This argument continues until we reach a stage where in our list we can no longer go one more step down in one of the m 's to make new states.

Let j_b be the smaller of j_a and j_b (if they are equal take either one); then only $2j_b$ values of J are required—going in integer steps from $j_a + j_b$ down to $j_a - j_b$. That is, when two objects of spin j_a and j_b are combined, the system can have a total angular momentum J equal to any one of the values

$$J = \begin{cases} j_a + j_b \\ j_a + j_b - 1 \\ j_a + j_b - 2 \\ \vdots \\ |j_a - j_b|. \end{cases} \quad (18.58)$$

(By writing $|j_a - j_b|$ instead of $j_a - j_b$ we can avoid the extra admonition that $j_a \geq j_b$.)

For *each* of these J values there are the $2J + 1$ states of different M -values—with M going from $+J$ to $-J$. Each of these is formed from linear combinations of the original states $|a, m_a; b, m_b\rangle$ with appropriate factors—the Clebsch-Gordon

coefficients for each particular term. We can consider that these coefficients give the “amount” of the state $|j_a, m_a; j_b, m_b\rangle$ which appears in the state $|J, M\rangle$. So each of the Clebsch-Gordon coefficients has, if you wish, *six* indices identifying its position in the formulas like those of Tables 18-3 and 18-6. That is, calling these coefficients $C(J, M; j_a, m_a; j_b, m_b)$, we could express the equality of the second line of Table 18-6 by writing

$$C(\frac{3}{2}, +\frac{1}{2}; \frac{1}{2}, +\frac{1}{2}; 1, 0) = \sqrt{2/3},$$

$$C(\frac{3}{2}, +\frac{1}{2}; \frac{1}{2}, -\frac{1}{2}; 1, +1) = \sqrt{1/3}.$$

We will not calculate here the coefficients for any other special cases.† You can, however, find tables in many books. You might wish to try another special case for yourself. The next one to do would be the composition of two spin-one particles. We give just the final result in Table 18-7.

These laws of the composition of angular momenta are very important in particle physics—where they have innumerable applications. Unfortunately, we have no time to look at more examples here.

Table 18-7

Composition of two spin-one particles ($j_a = 1, j_b = 1$)

$ J = 2, M = +2\rangle = a, +1; b, +1\rangle$
$ J = 2, M = +1\rangle = \frac{1}{\sqrt{2}} a, +1; b, 0\rangle + \frac{1}{\sqrt{2}} a, 0; b, +1\rangle$
$ J = 2, M = 0\rangle = \frac{1}{\sqrt{6}} a, +1; b, -1\rangle + \frac{1}{\sqrt{6}} a, -1; b, +1\rangle + \frac{2}{\sqrt{6}} a, 0; b, 0\rangle$
$ J = 2, M = -1\rangle = \frac{1}{\sqrt{2}} a, 0; b, -1\rangle + \frac{1}{\sqrt{2}} a, -1; b, 0\rangle$
$ J = 2, M = -2\rangle = a, -1; b, -1\rangle$
$ J = 1, M = +1\rangle = \frac{1}{\sqrt{2}} a, +1; b, 0\rangle - \frac{1}{\sqrt{2}} a, 0; b, +1\rangle$
$ J = 1, M = 0\rangle = \frac{1}{\sqrt{2}} a, +1; b, -1\rangle - \frac{1}{\sqrt{2}} a, -1; b, +1\rangle$
$ J = 1, M = -1\rangle = \frac{1}{\sqrt{2}} a, 0; b, -1\rangle - \frac{1}{\sqrt{2}} a, -1; b, 0\rangle$
$ J = 0, M = 0\rangle = \frac{1}{\sqrt{3}} \{ a, +1; b, -1\rangle + a, -1; b, +1\rangle - a, 0; b, 0\rangle \}$

Added Note 1: Derivation of the rotation matrix‡

For those who would like to see the details, we work out here the general rotation matrix for a system with spin (total angular momentum) j . It is really not very important to work out the general case; once you have the idea, you can find the general results in tables in many books. On the other hand, after coming this far you might like to see that you can indeed understand even the very complicated formulas of quantum mechanics, such as Eq. (18.35), that come into the description of angular momentum.

† A large part of the work is done now that we have the general rotation matrix Eq. (18.35).

‡ The material of this appendix was originally included in the body of the lecture. We now feel that it is unnecessary to include such a detailed treatment of the general case.

We extend the arguments of Section 18-4 to a system with spin j , which we consider to be made up of $2j$ spin one-half objects. The state with $m = j$ would be $|++++\cdots+\rangle$ (with j plus signs). For $m = j - 1$, there will be $2j$ terms like $|++++\cdots+-\rangle, |+++ \cdots + -\rangle$, and so on. Let's consider the general case in which there are r plusses and s minuses—with $r + s = 2j$. Under a rotation about the z -axis each of the r plusses will contribute $e^{+i\phi/2}$. The result is a phase change of $i(r/2 - s/2)\phi$. You see that

$$m = \frac{r - s}{2}. \quad (18.59)$$

Just as for $J = \frac{3}{2}$, each state of definite m must be the linear combination with plus signs of all the states with the same r and s —that is, states corresponding to every possible arrangement which has r plusses and s minuses. We assume that you can figure out that there are $(r + s)!/r!s!$ such arrangements. To normalize each state, we should divide the sum by the square root of this number. We can write

$$\left[\frac{(r + s)!}{r!s!} \right]^{-1/2} \{ \underbrace{++++\cdots+}_r \underbrace{---\cdots-}_s \} + (\text{all rearrangements of order}) = |j, m\rangle \quad (18.60)$$

with

$$j = \frac{r + s}{2}, \quad m = \frac{r - s}{2}. \quad (18.61)$$

It will help our work if we now go to still another notation. Once we have defined the states by Eq. (18.60), the two numbers r and s define a state just as well as j and m . It will help us keep track of things if we write

$$|j, m\rangle = |r_s\rangle, \quad (18.62)$$

where, using the equalities of (18.67)

$$r = j + m, \quad s = j - m.$$

Next, we would like to write Eq. (18.60) with a new *special notation* as

$$|j, m\rangle = |r_s\rangle = \left[\frac{(r + s)!}{r!s!} \right]^{+1/2} \{ |+\rangle^r |-\rangle^s \}_{\text{perm}}. \quad (18.63)$$

Note that we have changed the exponent of the factor in front to *plus* $\frac{1}{2}$. We do that because there are just $N = (r + s)!/r!s!$ terms inside the curly brackets. Comparing (18.63) with (18.60) it is clear that

$$\{ |+\rangle^r |-\rangle^s \}_{\text{perm}}$$

is just a shorthand way of writing

$$\frac{\{ |++++\cdots- \rangle + \text{all rearrangements} \}}{N},$$

where N is the number of different terms in the bracket. The reason that this notation is convenient is that each time we make a rotation, all of the plus signs contribute the same factor, so we get this factor to the r th power. Similarly, all together the s minus terms contribute a factor to the s th power no matter what the sequence of the terms is.

Now suppose we rotate our system by the angle θ about the y -axis. What we want is $R_y(\theta) |r_s\rangle$. When $R_y(\theta)$ operates on each $|+\rangle$ it gives

$$R_y(\theta) |+\rangle = |+\rangle C + |-\rangle S, \quad (18.64)$$

where $C = \cos \theta/2$ and $S = \sin \theta/2$. When $R_y(\theta)$ operates on each $|-\rangle$ it gives

$$R_y(\theta) |-\rangle = |-\rangle C - |+\rangle S.$$

So what we want is

$$\begin{aligned}
R_y(\theta) | \begin{smallmatrix} r \\ s \end{smallmatrix} \rangle &= \left[\frac{(r+s)!}{r!s!} \right]^{1/2} R_y(\theta) \{ |+\rangle^r |-\rangle^s \}_{\text{perm}} \\
&= \left[\frac{(r+s)!}{r!s!} \right]^{1/2} \{ (R_y(\theta) |+\rangle)^r (R_y(\theta) |-\rangle)^s \}_{\text{perm}} \\
&= \left[\frac{(r+s)!}{r!s!} \right]^{1/2} \{ |+\rangle C + |-\rangle S \}^r \{ |-\rangle C - |+\rangle S \}^s \}_{\text{perm}}. \quad (18.65)
\end{aligned}$$

Now each binomial has to be expanded out to its appropriate power and the two expressions multiplied together. There will be terms with $|+\rangle$ to all powers from zero to $(r+s)$. Let's look at all of the terms which have $|+\rangle$ to the r' power. They will appear always multiplied with $|-\rangle$ to the s' power, where $s' = 2j - r'$. Suppose we collect all such terms. For each permutation they will have some numerical coefficient involving the factors of the binomial expansion as well as the factors C and S . Suppose we call that factor $A_{r'}$. Then Eq. (18.65) will look like

$$R_y(\theta) | \begin{smallmatrix} r \\ s \end{smallmatrix} \rangle = \sum_{r'=0}^{r+s} \{ A_{r'} |+\rangle^{r'} |-\rangle^{s'} \}_{\text{perm}}. \quad (18.66)$$

Now let's say that we divide $A_{r'}$ by the factor $[(r'+s')!/r's'!]^{1/2}$ and call the quotient $B_{r'}$. Equation (18.66) is then equivalent to

$$R_y(\theta) | \begin{smallmatrix} r \\ s \end{smallmatrix} \rangle = \sum_{r'=0}^{r+s} B_{r'} \left[\frac{r'+s'}{r's'} \right]^{1/2} \{ |+\rangle^{r'} |-\rangle^{s'} \}_{\text{perm}}. \quad (18.67)$$

(We could just say that this equation defines $B_{r'}$ by the requirement that (18.67) gives the same expression that appears in (18.65).)

With this definition of $B_{r'}$ the remaining factors on the right-hand side of Eq. (18.67) are just the states $| \begin{smallmatrix} r' \\ s' \end{smallmatrix} \rangle$. So we have that

$$R_y(\theta) | \begin{smallmatrix} r \\ s \end{smallmatrix} \rangle = \sum_{r'=0}^{r+s} B_{r'} | \begin{smallmatrix} r' \\ s' \end{smallmatrix} \rangle, \quad (18.68)$$

with s' always equal to $r+s-r'$. This means, of course, that the coefficients $B_{r'}$ are just the matrix elements we want, namely

$$\langle \begin{smallmatrix} r' \\ s' \end{smallmatrix} | R_y(\theta) | \begin{smallmatrix} r \\ s \end{smallmatrix} \rangle = B_{r'}. \quad (18.69)$$

Now we just have to push through the algebra to find the various $B_{r'}$. Comparing (18.39) with (18.37)—and remembering that $r'+s' = r+s$ —we see that $B_{r'}$ is just the coefficient of $a^{r'}b^{s'}$ in the following expression:

$$\left(\frac{r'!s'!}{r!s!} \right)^{1/2} (aC + bS)^r (bC - aS)^s. \quad (18.70)$$

It is now only a dirty job to make the expansions by the binomial theorem, and collect the terms with the given power of a and b . If you work it all out, you find that the coefficient of $a^{r'}b^{s'}$ in (18.70) is

$$\left[\frac{r'!s'!}{r!s!} \right]^{1/2} \sum_k (-1)^k S^{r-r'+2k} C^{s+r'-2k} \cdot \frac{r!}{(r-r'+k)!(r'-k)!} \cdot \frac{s!}{(s-k)!k!}. \quad (18.71)$$

The sum is to be taken over all integers k which give terms of zero or greater in the factorials. This expression is then the matrix element we wanted.

Finally, we can return to our original notation in terms of j , m , and m' using

$$r = j + m, \quad r' = j + m', \quad s = j - m, \quad s' = j - m'.$$

Making these substitutions, we get Eq. (18.34) in Section 18-4.

Added Note 2: Conservation of parity in photon emission

In Section 1 of this chapter we considered the emission of light by an atom that goes from an excited state of spin 1 to a ground state of spin 0. If the excited state has its spin up ($m = +1$), it can emit a RHC photon along the $+z$ -axis or a LHC photon along the $-z$ -axis. Let's call these two states of the photon $|R_{up}\rangle$ and $|L_{dn}\rangle$. Neither of these states has a definite parity. Letting \hat{P} be the parity operator, $\hat{P}|R_{up}\rangle = |L_{dn}\rangle$ and $\hat{P}|L_{dn}\rangle = |R_{up}\rangle$.

What about our earlier proof that an atom in a state of definite energy must have a definite parity, and our statement that parity is conserved in atomic processes? Shouldn't the final state in this problem (the state after the emission of a photon) have a definite parity? It *does* if we consider the *complete* final state which contains amplitudes for the emission photons into all sorts of angles. In Section 1 we chose to consider only a part of the complete final state.

If we wish we can look only at final states that do have a definite parity. For example, consider a final state $|\psi_F\rangle$ which has some amplitude α to be a RHC photon going along $+z$ and some amplitude β to be a LHC photon going along $-z$. We can write

$$|\psi_F\rangle = \alpha |R_{up}\rangle + \beta |L_{dn}\rangle. \quad (18.72)$$

The parity operation on this state gives

$$\hat{P}|\psi_F\rangle = \alpha |L_{dn}\rangle + \beta |R_{up}\rangle. \quad (18.73)$$

This state will be $\pm |\psi_F\rangle$ if $\beta = \alpha$ or if $\beta = -\alpha$. So a final state of even parity is

$$|\psi_F^+\rangle = \alpha \{|R_{up}\rangle + |L_{dn}\rangle\}, \quad (18.74)$$

and a state of odd parity is

$$|\psi_F^-\rangle = \alpha \{|R_{up}\rangle - |L_{dn}\rangle\}. \quad (18.75)$$

Next, we wish to consider the decay of an excited state of odd parity to a ground state of even parity. If parity is to be conserved, the final state of the photon must have odd parity. It must be the state in (18.75). If the amplitude to find $|R_{up}\rangle$ is α , the amplitude to find $|L_{dn}\rangle$ is $-\alpha$.

Now notice what happens when we perform a rotation of 180° about the y -axis. The initial excited state of the atom becomes an $m = -1$ state (with no change in sign, according to Table 17-2). And the rotation of the final state gives

$$R_y(180^\circ)|\psi_F^-\rangle = \alpha \{|R_{dn}\rangle - |L_{up}\rangle\}. \quad (18.76)$$

Comparing this equation with (18.75), you see that for the assumed parity of the final state, the amplitude to get a LHC photon along $+z$ from the $m = -1$ initial state is the negative of the amplitude to get a RHC photon from the $m = +1$ initial state. This agrees with the result we found in Section 1.

The Hydrogen Atom and The Periodic Table

19-1 Schrödinger's equation for the hydrogen atom

The most dramatic success in the history of the quantum mechanics was the understanding of the details of the spectra of some simple atoms and the understanding of the periodicities which are found in the table of chemical elements. In this chapter we will at last bring our quantum mechanics to the point of this important achievement, specifically to an understanding of the spectrum of the hydrogen atom. We will at the same time arrive at a qualitative explanation of the mysterious properties of the chemical elements. We will do this by studying in detail the behavior of the electron in a hydrogen atom—for the first time making a detailed calculation of a distribution-in-space according to the ideas we developed in Chapter 16.

For a complete description of the hydrogen atom we should describe the motions of both the proton and the electron. It is possible to do this in quantum mechanics in a way that is analogous to the classical idea of describing the motion of each particle relative to the center of gravity, but we will not do so. We will just discuss an approximation in which we consider the proton to be very heavy, so we can think of it as fixed at the center of the atom.

We will make another approximation by forgetting that the electron has a spin and should be described by relativistic laws of mechanics. Some small corrections to our treatment will be required since we will be using the nonrelativistic Schrödinger equation and will disregard magnetic effects. Small magnetic effects occur because from the electron's point-of-view the proton is a circulating charge which produces a magnetic field. In this field the electron will have a different energy with its spin up than with it down. The energy of the atom will be shifted a little bit from what we will calculate. We will ignore this small energy shift. Also we will imagine that the electron is just like a gyroscope moving around in space always keeping the same direction of spin. Since we will be considering a free atom in space the total angular momentum will be conserved. In our approximation we will assume that the angular momentum of the electron spin stays constant, so all the rest of the angular momentum of the atom—what is usually called “orbital” angular momentum—will also be conserved. To an excellent approximation the electron moves in the hydrogen atom like a particle without spin—the angular momentum of the motion is a constant.

With these approximations the amplitude to find the electron at different places in space can be represented by a function of position in space and time. We let $\psi(x, y, z, t)$ be the amplitude to find the electron somewhere at the time t . According to the quantum mechanics the rate of change of this amplitude with time is given by the Hamiltonian operator working on the same function. From Chapter 16,

$$i\hbar \frac{\partial \psi}{\partial t} = \hat{\mathcal{H}}\psi, \quad (19.1)$$

with

$$\hat{\mathcal{H}} = -\frac{\hbar^2}{2m} \nabla^2 + V(\mathbf{r}). \quad (19.2)$$

Here, m is the electron mass, and $V(\mathbf{r})$ is the potential energy of the electron in the

19-1 Schrödinger's equation for the hydrogen atom

19-2 Spherically symmetric solutions

19-3 States with an angular dependence

19-4 The general solution for hydrogen

19-5 The hydrogen wave functions

19-6 The periodic table

electrostatic field of the proton. Taking $V = 0$ at large distances from the proton we can write†

$$V = -\frac{e^2}{r}.$$

The wave function ψ must then satisfy the equation

$$i\hbar \frac{\partial \psi}{\partial t} = -\frac{\hbar^2}{2m} \nabla^2 \psi - \frac{e^2}{r} \psi. \quad (19.3)$$

We want to look for definite energy states, so we try to find solutions which have the form

$$\psi(r, t) = e^{-(i/\hbar)Et} \psi(r). \quad (19.4)$$

The function $\psi(r)$ must then be a solution of

$$-\frac{\hbar^2}{2m} \nabla^2 \psi = \left(E + \frac{e^2}{r}\right) \psi, \quad (19.5)$$

where E is some constant—the energy of the atom.

Since the potential energy term depends only on the radius, it turns out to be much more convenient to solve this equation in polar coordinates rather than rectangular coordinates. The Laplacian is defined in rectangular coordinates by

$$\nabla^2 = \frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} + \frac{\partial^2}{\partial z^2}.$$

We want to use instead the coordinates r, θ, ϕ shown in Fig. 19-1. These coordinates are related to x, y, z by

$$x = r \sin \theta \cos \phi; \quad y = r \sin \theta \sin \phi; \quad z = r \cos \theta.$$

It's a rather tedious mess to work through the algebra, but you can eventually show that for any function $f(r) = f(r, \theta, \phi)$,

$$\nabla^2 f(r, \theta, \phi) = \frac{1}{r} \frac{\partial^2}{\partial r^2} (rf) + \frac{1}{r^2} \left\{ \frac{1}{\sin \theta} \frac{\partial}{\partial \theta} \left(\sin \theta \frac{\partial f}{\partial \theta} \right) + \frac{1}{\sin^2 \theta} \frac{\partial^2 f}{\partial \phi^2} \right\}. \quad (19.6)$$

So in terms of the polar coordinates, the equation which is to be satisfied by $\psi(r, \theta, \phi)$ is

$$\frac{1}{r} \frac{\partial^2}{\partial r^2} (r\psi) + \frac{1}{r^2} \left\{ \frac{1}{\sin \theta} \frac{\partial}{\partial \theta} \left(\sin \theta \frac{\partial \psi}{\partial \theta} \right) + \frac{1}{\sin^2 \theta} \frac{\partial^2 \psi}{\partial \phi^2} \right\} = -\frac{2m}{\hbar^2} \left(E + \frac{e^2}{r} \right) \psi. \quad (19.7)$$

19-2 Spherically symmetric solutions

Let's first try to find some very simple function that satisfies the horrible equation in (19.7). Although the wave function ψ will, in general, depend on the angles θ and ϕ as well as on the radius r , we can see whether there might be a special situation in which ψ does *not* depend on the angles. For a wave function that doesn't depend on the angles, none of the amplitudes will change in any way if you rotate the coordinate system. That means that all of the components of the angular momentum are zero. Such a ψ must correspond to a state whose total angular momentum is zero. (Actually, it is only the orbital angular momentum which is zero because we still have the spin of the electron, but we are ignoring that part.) A state with zero orbital angular momentum is called by a special name. It is called an “s-state”—you can remember “s for spherically symmetric.”‡

† As usual, $e^2 = q_e^2/4\pi\epsilon_0$.

‡ Since these special names are part of the common vocabulary of atomic physics, you will just have to learn them. We will help out by putting them together in a short “dictionary” later in the chapter.

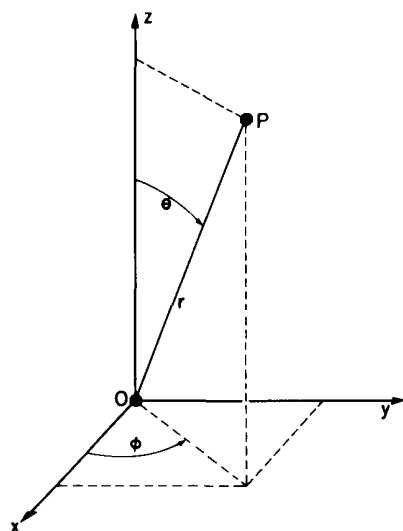


Fig. 19-1. The spherical polar coordinates r, θ, ϕ of the point P .

Now if ψ is not going to depend on θ and ϕ then the entire Laplacian contains only the first term and Eq. (19.7) becomes much simpler:

$$\frac{1}{r} \frac{d^2}{dr^2} (r\psi) = -\frac{2m}{\hbar^2} \left(E + \frac{e^2}{r} \right) \psi. \quad (19.8)$$

Before you start to work on solving an equation like this, it's a good idea to get rid of all excess constants like e^2 , m , and \hbar , by making some scale changes. Then the algebra will be easier. If we make the following substitutions:

$$r = \frac{\hbar^2}{me^2} \rho, \quad (19.9)$$

and

$$E = \frac{me^4}{2\hbar^2} \epsilon, \quad (19.10)$$

then Eq. (19.8) becomes (after multiplying through by ρ)

$$\frac{d^2(\rho\psi)}{d\rho^2} = -\left(\epsilon + \frac{2}{\rho} \right) \rho\psi. \quad (19.11)$$

These scale changes mean that we are measuring the distance r and energy E as multiples of “natural” atomic units. That is, $\rho = r/r_B$, where $r_B = \hbar^2/me^2$, is called the “Bohr radius” and is about 0.528 angstroms. Similarly, $\epsilon = E/E_R$, with $E_R = me^4/2\hbar^2$. This energy is called the “Rydberg” and is about 13.6 electron volts.

Since the product $\rho\psi$ appears on both sides, it is convenient to work with it rather than with ψ itself. Letting

$$\rho\psi = f, \quad (19.12)$$

we have the more simple-looking equation

$$\frac{d^2f}{d\rho^2} = -\left(\epsilon + \frac{2}{\rho} \right) f. \quad (19.13)$$

Now we have to find some function f which satisfies Eq. (19.13)—in other words, we just have to solve a differential equation. Unfortunately, there is no very useful, general method for solving any given differential equation. You just have to fiddle around. Our equation is not easy, but people have found that it can be solved by the following procedure. First, you replace f , which is some function of ρ , by a product of two functions

$$f(\rho) = e^{-\alpha\rho} g(\rho). \quad (19.14)$$

This just means that you are factoring $e^{-\alpha\rho}$ out of $f(\rho)$. You can certainly do that for any $f(\rho)$ at all. This just shifts our problem to finding the right function $g(\rho)$.

Sticking (19.14) into (19.13), we get the following equation for g :

$$\frac{d^2g}{d\rho^2} - 2\alpha \frac{dg}{d\rho} + \left(\frac{2}{\rho} + \epsilon + \alpha^2 \right) g = 0. \quad (19.15)$$

Since we are free to choose α , let's make

$$\alpha^2 = -\epsilon, \quad (19.16)$$

and get

$$\frac{d^2g}{d\rho^2} - 2\alpha \frac{dg}{d\rho} + \frac{2}{\rho} g = 0. \quad (19.17)$$

You may think we are no better off than we were at Eq. (19.13), but the happy thing about our new equation is that it can be solved easily in terms of a power series in ρ (It is possible, in principle, to solve (19.13) that way too, but it is

much harder.) We are saying that Eq. (19.17) can be satisfied by some $g(\rho)$ which can be written as a series,

$$g(\rho) = \sum_{k=1}^{\infty} a_k \rho^k, \quad (19.18)$$

in which the a_k are constant coefficients. Now all we have to do is find a suitable infinite set of coefficients! Let's check to see that such a solution will work. The first derivative of this $g(\rho)$ is

$$\frac{dg}{d\rho} = \sum_{k=1}^{\infty} a_k k \rho^{k-1},$$

and the second derivative is

$$\frac{d^2g}{d\rho^2} = \sum_{k=1}^{\infty} a_k k(k-1) \rho^{k-2}.$$

Using these expressions in (19.17) we have

$$\sum_{k=1}^{\infty} k(k-1)a_k \rho^{k-2} - \sum_{k=1}^{\infty} 2\alpha k a_k \rho^{k-1} + \sum_{k=1}^{\infty} 2a_k \rho^{k-1} = 0. \quad (19.19)$$

It's not obvious that we have succeeded; but we forge onward. It will all look better if we replace the first sum by an equivalent. Since the first term of the sum is zero, we can replace each k by $k+1$ without changing anything in the infinite series; with this change the first sum can equally well be written as

$$\sum_{k=1}^{\infty} (k+1)k a_{k+1} \rho^{k-1}.$$

Now we can put all the sums together to get

$$\sum_{k=1}^{\infty} [(k+1)k a_{k+1} - 2\alpha k a_k + 2a_k] \rho^{k-1} = 0. \quad (19.20)$$

This power series must vanish for all possible values of ρ . It can do that only if the coefficient of each power of ρ is separately zero. We will have a solution for the hydrogen atom if we can find a set a_k for which

$$(k+1)k a_{k+1} - 2(\alpha k - 1)a_k = 0 \quad (19.21)$$

for all $k > 1$. That is certainly easy to arrange. Pick any a_1 you like. Then generate all of the other coefficients from

$$a_{k+1} = \frac{2(\alpha k - 1)}{k(k+1)} a_k. \quad (19.22)$$

With this you will get a_2, a_3, a_4 , and so on, and each pair will certainly satisfy (19.21). We get a series for $g(\rho)$ which satisfies (19.17). With it we can make a ψ_1 that satisfies Schrödinger's equation. Notice that the solutions depend on the assumed energy (through α), but for each value of ϵ , there is a corresponding series.

We have a solution, but what does it represent physically? We can get an idea by seeing what happens far from the proton—for large values of ρ . Out there, the high-order terms of the series are the most important, so we should look at what happens for large k . When $k \gg 1$, Eq. (19.22) is approximately the same as

$$a_{k+1} = \frac{2\alpha}{k} a_k,$$

which means that

$$a_{k+1} \approx \frac{(2\alpha)^k}{k!}. \quad (19.23)$$

But these are just the coefficients of the series for $e^{+2\alpha\rho}$. The function of g is a rapidly increasing exponential. Even coupled with $e^{-\alpha\rho}$ to produce $f(\rho)$ —see

Eq. (19.14)—it still gives a solution for $f(\rho)$ which goes like $e^{\alpha\rho}$ for large ρ . We have found a mathematical solution but not a physical one. It represents a situation in which the electron is *least* likely to be near the proton! It is always more likely to be found at a very large radius ρ . A wave function for a *bound* electron must go to zero for large ρ .

We have to think whether there is some way to beat the game, and there is. Observe! If it just happened by luck that α were equal to $1/n$, where n is any integer, then Eq. (19.22) would make $a_{n+1} = 0$. All higher terms would also be zero. We wouldn't have an infinite series but a finite polynomial. Any polynomial increases more slowly than $e^{\alpha\rho}$, so the term $e^{-\alpha\rho}$ will eventually beat it down, and the function f will go to zero for large ρ . The only *bound-state* solutions are those for which $\alpha = 1/n$, with $n = 1, 2, 3, 4$, and so on.

Looking back to Eq. (19.16), we see that the bound-state solutions to the spherically symmetric wave equation can exist only when

$$-\epsilon = 1, \frac{1}{4}, \frac{1}{9}, \frac{1}{16}, \dots, \frac{1}{n^2}, \dots$$

The allowed energies are just these fractions times the Rydberg, $E_R = me^4/2\hbar^2$, or the energy of the n th energy level is

$$E_n = -E_R \frac{1}{n^2}. \quad (19.24)$$

There is, incidentally, nothing mysterious about negative numbers for the energy. The energies are negative because when we chose to write $V = -e^2/r$, we picked our zero point as the energy of an electron located far from the proton. When it is close to the proton, its energy is less, so somewhat below zero. The energy is lowest (most negative) for $n = 1$, and increases toward zero with increasing n .

Before the discovery of quantum mechanics, it was known from experimental studies of the spectrum of hydrogen that the energy levels could be described by Eq. (19.24), where E_R was found from the observations to be about 13.6 electron volts. Bohr then devised a model which gave the same equation and predicted that E_R should be $me^4/2\hbar^2$. But it was the first great success of the Schrödinger theory that it could reproduce this result from a basic equation of motion for the electron.

Now that we have solved our first atom, let's look at the nature of the solution we got. Pulling all the pieces together, each solution looks like this:

$$\psi_n = \frac{f_n(\rho)}{\rho} = \frac{e^{-\rho/n}}{\rho} g_n(\rho), \quad (19.25)$$

where

$$g_n(\rho) = \sum_{k=1}^n a_k \rho^k \quad (19.26)$$

and

$$a_{k+1} = \frac{2(k/n - 1)}{k(k+1)} a_k. \quad (19.27)$$

So long as we are mainly interested in the relative probabilities of finding the electron at various places we can pick any number we wish for a_1 . We may as well set $a_1 = 1$. (People often choose a_1 so that the wave function is "normalized," that is, so that the integrated probability of finding the electron anywhere in the atom is equal to 1. We have no need to do that just now.)

For the lowest energy state, $n = 1$, and

$$\psi_1(\rho) = e^{-\rho}. \quad (19.28)$$

For a hydrogen atom in its ground (lowest-energy) state, the amplitude to find the electron at any point drops off exponentially with the distance from the proton. It is most likely to be found right at the proton, and the characteristic spreading distance is about one unit in ρ , or about one Bohr radius, r_B .

Putting $n = 2$ gives the next higher level. The wave function for this state will have two terms. It is

$$\psi_2(\rho) = \left(1 - \frac{\rho}{2}\right) e^{-\rho/2}. \quad (19.29)$$

The wave function for the next level is

$$\psi_3(\rho) = \left(1 - \frac{2\rho}{3} + \frac{2}{27}\rho^2\right) e^{-\rho/3}. \quad (19.30)$$

The wave functions for these first three levels are plotted in Fig. 19-2. You can see the general trend. All of the wave functions approach zero rapidly for large ρ after oscillating a few times. In fact, the number of “bumps” is just equal to n —or, if you prefer, the number of zero-crossings of ψ_n is $n - 1$.

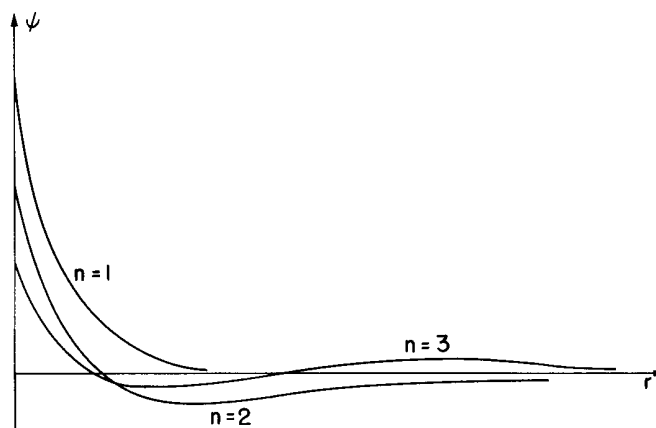


Fig. 19-2. The wave functions for the first three $l = 0$ states of the hydrogen atom. (The scales are chosen so that the total probabilities are equal.)

19-3 States with an angular dependence

In the states described by the $\psi_n(r)$ we have found that the probability amplitude for finding the electron is spherically symmetric—depending only on r , the distance from the proton. Such states have zero orbital angular momentum. We should now inquire about states which may have some angular dependences.

We could, if we wished, just investigate the strictly mathematical problem of finding the functions of r , θ , and ϕ which satisfy the differential equation (19.7)—putting in the additional physical conditions that the only acceptable functions are ones which go to zero for large r . You will find this done in many books. We are going to take a short cut by using the knowledge we already have about how amplitudes depend on angles in space.

The hydrogen atom in any particular state is a particle with a certain “spin” J —the quantum number of the total angular momentum. Part of this spin comes from the electron’s intrinsic spin, and part from the electron’s motion. Since each of these two components acts independently (to an excellent approximation) we will again ignore the spin part and think only about the “orbital” angular momentum. This orbital motion behaves, however, just like a spin. For example, if the orbital quantum number is l , the z -component of angular momentum can be $l, l - 1, l - 2, \dots, -l$. (We are, as usual, measuring in units of \hbar .) Also, all the rotation matrices and other properties we have worked out still apply (From now on we will *really* ignore the electron’s spin; when we speak of “angular momentum” we will mean only the orbital part.)

Since the potential V in which the electron moves depends only on r and not on θ or ϕ , the Hamiltonian is symmetric under all rotations. It follows that the angular momentum and all its components are conserved. (This is true for motion in *any* “central field”—one which depends only on r —so is not a special feature of the Coulomb e^2/r potential.)

Now let's think of some possible state of the electron; its internal angular structure will be characterized by the quantum number l . Depending on the "orientation" of the total angular momentum with respect to the z -axis, the z -component of angular momentum will be m , which is one of the $2l + 1$ possibilities between $+l$ and $-l$. Let's say $m = 1$. With what amplitude will the electron be found on the z -axis at some distance r ? Zero. An electron on the z -axis *cannot* have any orbital angular momentum around that axis. Alright, suppose m is zero, then there can be some nonzero amplitude to find the electron at each distance from the proton. We'll call this amplitude $F_l(r)$. It is the amplitude to find the electron at the distance r up along the z -axis, when the atom is in the state $|l, 0\rangle$, by which we mean orbital spin l and z -component $m = 0$.

If we know $F_l(r)$ everything is known. For any state $|l, m\rangle$, we know the amplitude $\psi_{l,m}(r)$ to find the electron *anywhere* in the atom. How? Watch. Suppose we have the atom in the state $|l, m\rangle$, what is the amplitude to find the electron at the angle θ, ϕ and the distance r from the origin? Put a new z -axis, say z' , at that angle (see Fig. 19-3), and ask what is the amplitude that the electron will be at the distance r along the new axis z' ? We know that it cannot be found along z' unless its z' -component of angular momentum, say m' , is zero. When m' is zero, however, the amplitude to find the electron along z' is $F_l(r)$. Therefore, the result is the product of two factors. The first is the amplitude that an atom in the state $|l, m\rangle$ along the z -axis will be in the state $|l, m' = 0\rangle$ with respect to the z' -axis. Multiply that amplitude by $F_l(r)$ and you have the amplitude $\psi_{l,m}(r)$ to find the electron at (r, θ, ϕ) with respect to the original axes.

Let's write it out. We have worked out earlier the transformation matrices for rotations. To go from the frame x, y, z to the frame x', y', z' of Fig. 19-3, we can rotate first around the z -axis by the angle ϕ , and then rotate about the *new* y -axis (y') by the angle θ . This combined rotation is the product

$$R_y(\theta)R_z(\phi).$$

The amplitude to find the state $l, m' = 0$ after the rotation is

$$\langle l, 0 | R_y(\phi)R_z(\phi) | l, m \rangle. \quad (19.31)$$

Our result, then, is

$$\psi_{l,m}(r) = \langle l, 0 | R_y(\theta)R_z(\phi) | l, m \rangle F_l(r). \quad (19.32)$$

The orbital motion can have only integral values of l . (If the electron can be found anywhere at $r \neq 0$, there is some amplitude to have $m = 0$ in that direction. And $m = 0$ states exist only for integral spins.) The rotation matrices for $l = 1$ are given in Table 17-2. For larger l you can use the general formulas we worked out in Chapter 18. The matrices for $R_z(\phi)$ and $R_y(\theta)$ appear separately, but you know how to combine them. For the general case you would start with the state $|l, m\rangle$ and operate with $R_z(\phi)$ to get the new state $R_z(\phi) | l, m \rangle$. Then you operate on this state with $R_y(\theta)$ to get the state $R_y(\theta)R_z(\phi) | l, m \rangle$ (which is just $e^{im\phi} | l, m \rangle$). Multiplying by $\langle l, 0 |$ gives the matrix element (19.31).

The matrix elements of the rotation operation are algebraic functions of θ and ϕ . The particular functions which appear in (19.31) also show up in many kinds of problems which involve waves in spherical geometries and so has been given a special name. Not everyone uses the same convention; but one of the most common ones is

$$\langle l, 0 | R_y(\theta)R_z(\phi) | l, m \rangle \equiv a Y_{l,m}(\theta, \phi). \quad (19.33)$$

The functions $Y_{l,m}(\theta, \phi)$ are called the *spherical harmonics*, and a is just a numerical factor which depends on the definition chosen for $Y_{l,m}$. For the usual definition,

$$a = \sqrt{\frac{4}{2l+1}}. \quad (19.34)$$

With this notation, the hydrogen wave functions can be written

$$\psi_{l,m}(r) = Y_{l,m}(\theta, \phi) F_l(r). \quad (19.35)$$

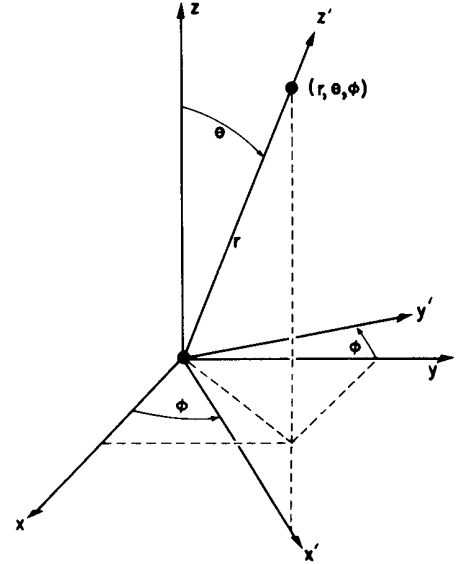
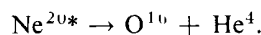


Fig. 19-3. The point (r, θ, ϕ) is on the z' -axis of the $x'y'z'$ coordinate frame.

The angle functions $Y_{l,m}(\theta, \phi)$ are important not only in many quantum-mechanical problems, but also in many areas of classical physics in which the ∇^2 operator appears, such as electromagnetism. As another example of their use in quantum mechanics, consider the disintegration of an excited state of Ne^{20} (such as we discussed in the last chapter) which decays by emitting an α -particle and going into O^{16} :



Suppose that the excited state has some spin l (necessarily an integer) and that the z -component of angular momentum is m . We might now ask the following: given l and m , what is the amplitude that we will find the α -particle going off in a direction which makes the angle θ with respect to the z -axis and the angle ϕ with respect to the xz -plane—as shown in Fig. 19-4.

To solve this problem we make, first, the following observation. A decay in which the α -particle goes straight up along z must come from a state with $m = 0$. This is so because both O^{16} and the α -particle have spin zero, and because their motion cannot have any angular momentum about the z -axis. Let's call this amplitude a (per unit solid angle). Then, to find the amplitude for a decay at the arbitrary angle of Fig. 19-4, all we need to know is what amplitude the given initial state has zero angular momentum about the decay direction. The amplitude for the decay at θ and ϕ is then a times the amplitude that a state $|l, m\rangle$ with respect to the z -axis will be in the state $|l, 0\rangle$ with respect to z' —the decay direction. This latter amplitude is just what we have written in (19.31). The probability to see the α -particle at θ, ϕ is

$$P(\theta, \phi) = a^2 |\langle l, 0 | R_y(\theta)R_z(\phi) | l, m \rangle|^2.$$

As an example, consider an initial state with $l = 1$ and various values of m . From Table 17-2 we know the necessary amplitudes. They are

$$\begin{aligned} \langle 1, 0 | R_y(\theta)R_z(\phi) | 1, +1 \rangle &= -\frac{1}{\sqrt{2}} \sin \theta e^{i\phi}, \\ \langle 1, 0 | R_y(\theta)R_z(\phi) | 1, 0 \rangle &= \cos \theta, \\ \langle 1, 0 | R_y(\theta)R_z(\phi) | 1, -1 \rangle &= -\frac{1}{\sqrt{2}} \sin \theta e^{-i\phi}. \end{aligned} \quad (19.36)$$

These are the three possible angular distribution amplitudes—depending on the m -value of the initial nucleus.

Amplitudes such as the ones in (19.36) appear so often and are sufficiently important that they are given *several* names. If the angular distribution amplitude is proportional to any one of the three functions or any linear combination of them, we say, "The system has an orbital angular momentum of one." Or we may say, "The Ne^{20*} emits a p -wave α -particle." Or we say, "The α -particle is emitted in an $l = 1$ state." Because there are so many ways of saying the same thing it is useful to have a dictionary. If you are going to understand what other physicists are talking about, you will just have to memorize the language. In Table 19-1 we give a dictionary of orbital angular momentum.

If the orbital angular momentum is zero, then there is no change when you rotate the coordinate system and there is no variation with angle—the "dependence" on angle is as a constant, say 1. This is also called an " s -state", and there is only one such state—as far as angular dependence is concerned. If the orbital angular momentum is 1, then the amplitude of the angular variation may be any one of the three functions given—depending on the value of m —or it may be a linear combination. These are called " p -states," and there are three of them. If the orbital angular momentum is 2 then there are the five functions shown. Any linear combination is called an " $l = 2$," or a " d -wave" amplitude. Now you can immediately guess what the next letter is—what should come after s, p, d ? Well, of course, f, g, h , and so on down the alphabet! The letters don't mean anything. (They did once mean something—they meant "sharp" lines, "principal" lines, "diffuse" lines and

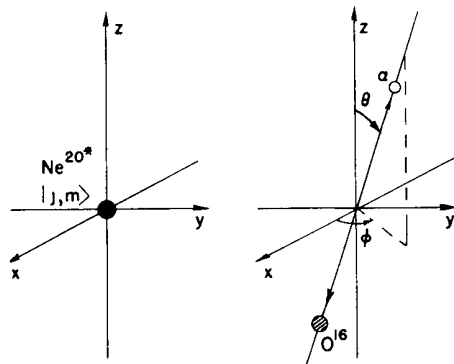


Fig. 19-4. The decay of an excited state of Ne^{20} .

Table 19-1

Dictionary of orbital angular momentum

 $(l = j = \text{an integer})$

Orbital angular momentum, l	z -component, m	Angular dependence of amplitudes	Name	Number of states	Orbital parity
0	0	1	s	1	+
1	$\begin{cases} +1 \\ 0 \\ -1 \end{cases}$	$\begin{cases} -\frac{1}{\sqrt{2}} \sin \theta e^{i\phi} \\ \cos \theta \\ \frac{1}{\sqrt{2}} \sin \theta e^{-i\phi} \end{cases}$	p	3	—
2	$\begin{cases} +2 \\ +1 \\ 0 \\ -1 \\ -2 \end{cases}$	$\begin{cases} \frac{\sqrt{6}}{4} \sin^2 \theta e^{2i\phi} \\ \frac{\sqrt{6}}{2} \sin \theta \cos \theta e^{i\phi} \\ \frac{1}{2}(3 \cos^2 \theta - 1) \\ -\frac{\sqrt{6}}{2} \sin \theta \cos \theta e^{-i\phi} \\ \frac{\sqrt{6}}{4} \sin^2 \theta e^{-2i\phi} \end{cases}$	d	5	+
3	$\left. \begin{matrix} \\ \\ \\ \end{matrix} \right\}$	$\left. \begin{matrix} l, 0 R_y(\theta) R_z(\phi) l, m \\ = Y_{l,m}(\theta, \phi) \\ = P_l^m(\cos \theta) e^{im\phi} \end{matrix} \right\}$	f	$2l + 1$	$(-1)^l$
4			g		
5			h		
\vdots			\vdots		

“fundamental” lines of the optical spectra of atoms. But those were in the days when people did not know where the lines came from. After f there were no special names, so we now just continue with g , h , and so on.)

The angular functions in the table go by several names—and are sometimes defined with slightly different conventions about the numerical factors that appear out in front. Sometimes they are called “spherical harmonics,” and written as $Y_{l,m}(\theta, \phi)$. Sometimes they are written $P_l^m(\cos \theta)e^{im\phi}$, and if $m = 0$, simply as $P_l(\cos \theta)$. The functions $P_l(\cos \theta)$ are called the “Legendre polynomials” in $\cos \theta$, and the functions $P_l^m(\cos \theta)$ are called the “associated Legendre functions.” You will find tables of these functions in many books.

Notice, incidentally, that all the functions for a given l have the property that they have the same parity—for odd l they change sign under an inversion and for even l they don’t change. So we can write that the parity of a state of orbital angular momentum l is $(-1)^l$.

As we have seen, these angular distributions may refer to a nuclear disintegration or some other process, or to the distribution of the amplitude to find an electron at some place in the hydrogen atom. For instance, if an electron is in a p -state ($l = 1$) the amplitude to find it can depend on the angle in many possible ways—but all are linear combinations of the three functions for $l = 1$ in Table 19-1. Let’s take the case $\cos \theta$. That’s interesting. That means that the amplitude is positive, say, in the upper part ($\theta < \pi/2$), is negative in the lower part ($\theta > \pi/2$), and is zero when θ is 90° . Squaring this amplitude we see that the probability of finding the electron varies with θ as shown in Fig. 19-5—and is independent of ϕ . This angular distribution is responsible for the fact that in molecular binding the attraction of an electron in an $l = 1$ state for another atom depends on direction—it is the origin of the directed valences of chemical attraction.

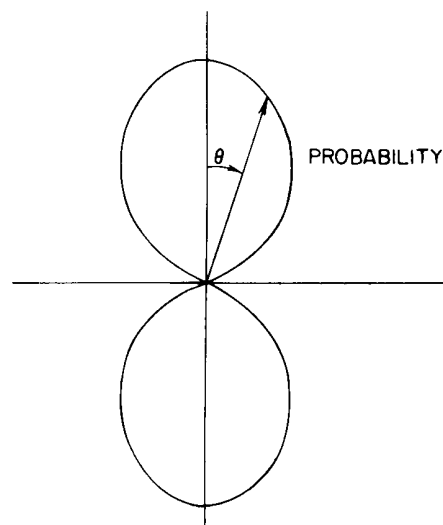


Fig. 19-5. A polar graph of $\cos^2 \theta$, which is the relative probability of finding an electron at various angles from the z -axis (for a given r) in an atomic state with $l = 1$ and $m = 0$.

19-4 The general solution for hydrogen

In Eq. (19.35) we have written the wave functions for the hydrogen atom as

$$\psi_{l,m}(r) = Y_{l,m}(\theta, \phi) F_l(r). \quad (19.37)$$

These wave functions must be solutions of the differential equation (19.7). Let's see what that means. Put (19.37) into (19.7); you get

$$\begin{aligned} \frac{Y_{l,m}}{r} \frac{\partial^2}{\partial r^2} (r F_l) + \frac{F_l}{r^2 \sin \theta} \frac{\partial}{\partial \theta} \left(\sin \theta \frac{\partial Y_{l,m}}{\partial \theta} \right) + \frac{F_l}{r^2 \sin^2 \theta} \frac{\partial^2 Y_{l,m}}{\partial \phi^2} \\ = - \frac{2m}{\hbar^2} \left(E + \frac{e^2}{r} \right) Y_{l,m} F_l. \end{aligned} \quad (19.38)$$

Now multiply through by r^2/F_l and rearrange terms. The result is

$$\begin{aligned} \frac{1}{\sin \theta} \frac{\partial}{\partial \theta} \left(\sin \theta \frac{\partial Y_{l,m}}{\partial \theta} \right) + \frac{1}{\sin^2 \theta} \frac{\partial^2 Y_{l,m}}{\partial \phi^2} \\ = - \left[\frac{r^2}{F_l} \left\{ \frac{1}{r} \frac{d^2}{dr^2} (r F_l) + \frac{2m}{\hbar^2} \left(E + \frac{e^2}{r} \right) \right\} \right] Y_{l,m}. \end{aligned} \quad (19.39)$$

The left-hand side of this equation depends on θ and ϕ , but *not* on r . No matter what value we choose for r , the left side doesn't change. *This must also be true for the right-hand side.* Although the quantity in the square brackets has r 's all over the place, the whole quantity cannot depend on r , otherwise we wouldn't have an equation good for all r . As you can see, the bracket also does not depend on θ or ϕ . It must be some constant. Its value may well depend on the l -value of the state we are studying, since the function F_l must be the one appropriate to that state; we'll call the constant K_l . Equation (19.35) is therefore equivalent to *two* equations:

$$\frac{1}{\sin \theta} \frac{\partial}{\partial \theta} \left(\sin \theta \frac{\partial Y_{l,m}}{\partial \theta} \right) + \frac{1}{\sin^2 \theta} \frac{\partial^2 Y_{l,m}}{\partial \phi^2} = -K_l Y_{l,m}, \quad (19.40)$$

$$\frac{1}{r} \frac{\partial^2}{\partial r^2} (r F_l) + \frac{2m}{\hbar^2} \left(E + \frac{e^2}{r} \right) F_l = K_l \frac{F_l}{r^2}. \quad (19.41)$$

Now look at what we've done. For any state described by l and m , we know the functions $Y_{l,m}$; we can use Eq. (19.40) to determine the constant K_l . Putting K_l into Eq. (19.41) we have a differential equation for the function $F_l(r)$. If we can solve that equation for $F_l(r)$, we have all of the pieces to put into (19.37) to give $\psi(r)$.

What is K_l ? First, notice that it must be the same for all m (which go with a particular l), so we can pick any m we want for $Y_{l,m}$ and plug it into (19.40) to solve for K_l . Perhaps the easiest one to use is $Y_{l,l}$. From Eq. (18.24),

$$R_z(\phi) |l, l\rangle = e^{il\phi} |l, l\rangle. \quad (19.42)$$

The matrix element for $R_y(\theta)$ is also quite simple:

$$\langle l, 0 | R_y(\theta) | l, l \rangle = b (\sin \theta)^l, \quad (19.43)$$

where b is some number.[†] Combining the two, we obtain

$$Y_{l,l} \propto e^{il\phi} \sin^l \theta. \quad (19.44)$$

[†] You can with some work show that this comes out of Eq. (18.35), but it is also easy to work out from first principles following the ideas of Section 18-4. A state $|l, l\rangle$ can be made out of $2l$ spin one-half particles all with spins up; while the state $|l, 0\rangle$ would have l up and l down. Under the rotation the amplitude that an up-spin remains up is $\cos \theta/2$, and that an up-spin goes down is $\sin \theta/2$. We are asking for the amplitude that l up-spins stay up, while the other l up-spins go down. The amplitude for that is $(\cos \theta/2 \sin \theta/2)^l$ which is the same as $\sin^l \theta$.

Putting this function into (19.40) gives

$$K_l = l(l + 1). \quad (19.45)$$

Now that we have determined K_l , Eq. (19.41) tells us about the radial function $F_l(r)$. It is, of course, just the Schrodinger equation with the angular part replaced by its equivalent $K_l F_l/r^2$. Let's rewrite (19.41) in the form we had in Eq. (19.8), as follows:

$$\frac{1}{r} \frac{d^2}{dr^2} (r F_l) = -\frac{2m}{\hbar^2} \left\{ E + \frac{e^2}{r} - \frac{l(l+1)\hbar^2}{2mr^2} \right\} F_l. \quad (19.46)$$

A mysterious term has been added to the potential energy. Although we got this term by some mathematical shenanigan, it has a simple physical origin. We can give you an idea about where it comes from in terms of a semi-classical argument. Then perhaps you will not find it quite so mysterious.

Think of a classical particle moving around some center of force. The total energy is conserved and is the sum of the potential and kinetic energies

$$U = V(r) + \frac{1}{2}mv^2 = \text{constant}.$$

In general, v can be resolved into a radial component v_r and a tangential component $r\dot{\theta}$; then

$$v^2 = v_r^2 + (r\dot{\theta})^2.$$

Now the angular momentum $mr^2\dot{\theta}$ is also conserved; say it is equal to L . We can then write

$$mr^2\dot{\theta} = L, \quad \text{or} \quad r\dot{\theta} = \frac{L}{mr},$$

and the energy is

$$U = \frac{1}{2}mv_r^2 + V(r) + \frac{L^2}{2mr^2}.$$

If there were no angular momentum we would have just the first two terms. Adding the angular momentum L does to the energy just what adding a term $L^2/2mr^2$ to the potential energy would do. But this is almost exactly the extra term in (19.46). The only difference is that $l(l+1)\hbar$ appears for the angular momentum instead of $L^2\hbar^2$ as we might expect. But we have seen before (for example, Volume II, Section 34-7)† that this is just the substitution that is usually required to make a quasi-classical argument agree with a correct quantum-mechanical calculation. We can, then, understand the new term as a “pseudo-potential” which gives the “centrifugal force” term that appears in the equations of radial motion for a rotating system. (See the discussion of “pseudo-forces” in Volume I, Section 12-5.)

We are now ready to solve Eq. (19.46) for $F_l(r)$. It is very much like Eq. (19.8), so the same technique will work again. Everything goes as before until you get to Eq. (19.19) which will have the additional term

$$-l(l+1) \sum_{k=1}^{\infty} a_k \rho^{k-2}. \quad (19.47)$$

This term can also be written as

$$-l(l+1) \left\{ \frac{a_1}{\rho} - \sum_{k=1}^{\infty} a_{k+1} \rho^{k-1} \right\}. \quad (19.48)$$

(We have taken out the first term and then shifted the running index k down by 1.) Instead of Eq. (19.20) we have

$$\sum_{k=1}^{\infty} [\{k(k+1) - l(l+1)\} a_{k+1} - 2(\alpha k - 1) a_k] \rho^{k-1} - \frac{l(l+1)a_1}{\rho} = 0. \quad (19.49)$$

† See Appendix to this volume.

There is only one term in ρ^{-1} , so it must be zero. The coefficient a_1 must be zero (unless $l = 0$ and we have our previous solution). Each of the other terms is made zero by having the square bracket come out zero for every k . This condition replaces Eq. (19.21) by

$$a_{k+1} = \frac{2(\alpha k - 1)}{k(k+1) - l(l+1)} a_k. \quad (19.50)$$

This is the only significant change from the spherically symmetric case.

As before the series must terminate if we are to have solutions which can represent bound electrons. The series will end at $k = n$ if $\alpha n = 1$. We get again the same condition on α , that it must be equal to $1/n$, where n is some integer. However, Eq. (19.50) also gives a new restriction. The index k cannot be equal to l , the denominator becomes zero and a_{l+1} is infinite. That is, since $a_1 = 0$, Eq. (19.50) implies that all successive a_k are zero until we get to a_{l+1} , which can be nonzero. This means that k must start at $l + 1$ and end at n .

Our final result is that for any l there are many possible solutions which we can call $F_{n,l}$ where $n \geq l + 1$. Each solution has the energy

$$E_n = -\frac{me^4}{2\hbar^2} \left(\frac{1}{n^2} \right). \quad (19.51)$$

The wave function for the state of this energy with the angular quantum numbers l and m is

$$\psi_{n,l,m} = Y_{l,m}(\theta, \phi) F_{n,l}(\rho), \quad (19.52)$$

with

$$\rho F_{n,l}(\rho) = e^{-\alpha\rho} \sum_{k=l+1}^n a_k \rho^k. \quad (19.53)$$

The coefficients a_k are obtained from (19.50). We have, finally, a complete description of the states of a hydrogen atom.

19-5 The hydrogen wave functions

Let's review what we have discovered. The states which satisfy Schrödinger's equation for an electron in a Coulomb field are characterized by three quantum numbers n, l, m , all integers. The angular distribution of the electron amplitude can have only certain forms which we call $Y_{l,m}$. They are labeled by l , the *quantum number of total angular momentum*, and m , the "*magnetic*" quantum number, which can range from $-l$ to $+l$. For each angular configuration, various possible radial distributions $F_{n,l}(r)$ of the electron amplitude are possible; they are labeled by the *principle quantum number* n —which can range from $l + 1$ to ∞ . The energy of the state depends only on n , and increases with increasing n .

The lowest energy, or ground, state is an s -state. It has $l = 0$, $n = 1$, and $m = 0$. It is a "nondegenerate" state—there is only one with this energy, and its wave function is spherically symmetric. The amplitude to find the electron is a maximum at the center, and falls off monotonically with increasing distance from the center. We can visualize the electron amplitude as a blob as shown in Fig. 19-6(a).

There are other s -states with higher energies, for $n = 2, 3, 4, \dots$. For each energy there is only one version ($m = 0$), and they are all spherically symmetric. These states have amplitudes which alternate in sign one or more times with increasing r . There are $n - 1$ spherical nodal surfaces—the places where ψ goes through zero. The $2s$ -state ($l = 0, n = 2$), for example, will look as sketched in Fig. 19-6(b). (The dark areas indicate regions where the amplitude is large, and the plus and minus signs indicate the relative phases of the amplitude.) The energy levels of the s -states are shown in the first column of Fig. 19-7.

Then there are the p -states—with $l = 1$. For each n , which must be 2 or greater, there are three states of the same energy, one each for $m = +1$, $m = 0$, and $m = -1$. The energy levels are as shown in Fig. 19-7. The angular dependences of these states are given in Table 19-1. For instance, for $m = 0$, if the

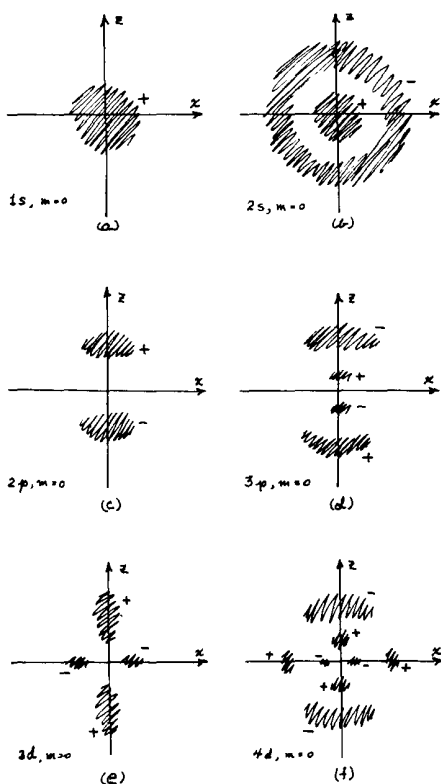


Fig. 19-6. Rough sketches showing the general nature of some of the hydrogen wave functions. The shaded regions show where the amplitudes are large. The plus and minus signs show the relative sign of the amplitude in each region.

amplitude is positive for θ near zero, it will be negative for θ near 180° . There is a nodal plane coincident with the xy -plane. For $n > 2$ there are also spherical nodes. The $n = 2, m = 0$ amplitude is sketched in Fig. 19-6(c), and the $n = 3, m = 0$ wave function is sketched in Fig. 19-6(d).

You might think that since m represents a kind of "orientation" in space, there should be similar distributions with the peaks of amplitude along the x -axis or along the y -axis. Are these perhaps the $m = +1$ and $m = -1$ states? No. But since we have three states with equal energies, any linear combinations of the three will also be stationary states of the same energy. It turns out that the "x"-state—which corresponds to the "z"-state, or $m = 0$ state, of Fig. 19-6(c)—is a linear combination of the $m = +1$ and $m = -1$ states. The corresponding "y"-state is another combination. Specifically, we mean that

$$\begin{aligned} \text{"z"} &= |1, 0\rangle, \\ \text{"x"} &= \frac{|1, +1\rangle + |1, -1\rangle}{\sqrt{2}}, \\ \text{"y"} &= \frac{|1, +1\rangle - |1, -1\rangle}{i\sqrt{2}}. \end{aligned}$$

These states all look the same when referred to their particular axes.

The d -states ($l = 2$) have five possible values of m for each energy, the lowest energy has $n = 3$. The levels go as shown in Fig. 19-7. The angular dependences get more complicated. For instance the $m = 0$ states have two conical nodes, so the wave function reverses phase from $+$, to $-$, to $+$ as you go around from the north pole to the south pole. The rough form of the amplitude is sketched in (e) and (f) of Fig. 19-6 for the $m = 0$ states with $n = 3$ and $n = 4$. Again, the larger n 's have spherical nodes.

We will not try to describe any more of the possible states. You will find the hydrogen wave functions described in more detail in many books. Two good references are L. Pauling and E. B. Wilson, *Introduction to Quantum Mechanics*, McGraw-Hill (1935); and R. B. Leighton, *Principles of Modern Physics*, McGraw-Hill (1959). You will find in them graphs of some of the functions and pictorial representations of many states.

We would like to mention one particular feature of the wave functions for higher l : for $l > 0$ the amplitudes are zero at the center. That is not surprising, since it's hard for an electron to have angular momentum when its radius arm is very small. For this reason, the higher the l , the more the amplitudes are "pushed away" from the center. If you look at the way the radial functions $F(r)$ vary for small r , you find from (19.53) that

$$F_{n,l}(r) \approx r^l.$$

Such a dependence on r means that for larger l 's you have to go farther from $r = 0$ before you get an appreciable amplitude. This behavior is, incidentally, determined by the centrifugal force term in the radial equation, so the same thing will apply for any potential that varies slower than $1/r^2$ for small r —which most atomic potentials do.

19-6 The periodic table

We would like now to apply the theory of the hydrogen atom in an approximate way to get some understanding of the chemist's periodic table of the elements. For an element with atomic number Z there are Z electrons held together by the electric attraction of the nucleus but with mutual repulsion of the electrons. To get an exact solution we would have to solve Schrödinger's equation for Z electrons in a Coulomb field. For helium the equation is

$$-\frac{\hbar}{i} \frac{\partial \psi}{\partial t} = -\frac{\hbar^2}{2m} (\nabla_1^2 \psi + \nabla_2^2 \psi) + \left(-\frac{2e^2}{r_1} - \frac{2e^2}{r_2} + \frac{e^2}{r_{12}} \right) \psi,$$

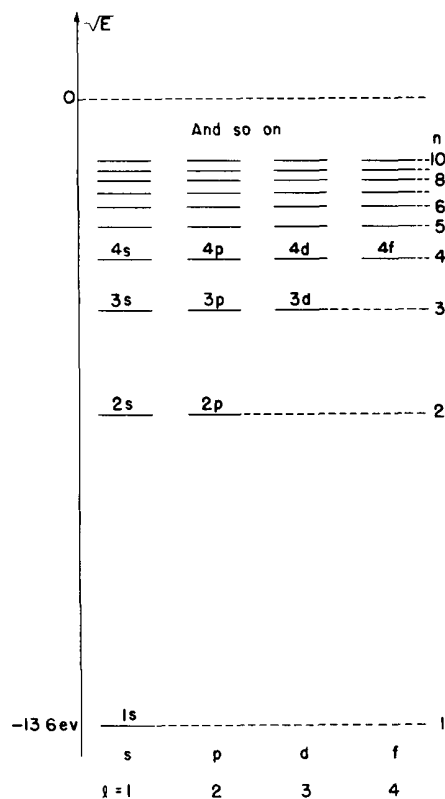


Fig. 19-7. The energy level diagram for hydrogen.

where ∇_1^2 is a Laplacian which operates on \mathbf{r}_1 , the coordinate of one electron; ∇_2^2 operates on \mathbf{r}_2 ; and $r_{12} = |\mathbf{r}_1 - \mathbf{r}_2|$. (We are again neglecting the spin of the electrons.) To find the stationary states and energy levels we would have to find solutions of the form

$$\psi = f(\mathbf{r}_1, \mathbf{r}_2)e^{-(i/\hbar)Et}.$$

The geometrical dependence is contained in f , which is a function of six variables—the simultaneous positions of the two electrons. No one has found an analytic solution, although solutions for the lowest energy states have been obtained by numerical methods.

With 3, 4, or 5 electrons it is hopeless to try to obtain exact solutions, and it is going too far to say that quantum mechanics has given a precise understanding of the periodic table. It is possible, however, even with a sloppy approximation—and some fixing—to understand, at least qualitatively, many chemical properties which show up in the periodic table.

The chemical properties of atoms are determined primarily by their lowest energy states. We can use the following approximate theory to find these states and their energies. First, we neglect the electron spin, *except* that we adopt the exclusion principle and say that any particular electronic state can be occupied by only one electron. This means that any particular orbital configuration can have up to *two* electrons—one with spin up, the other with spin down. Next we disregard the *details* of the interactions between the electrons in our first approximation, and say that each electron moves in a *central field* which is the combined field of the nucleus and all the other electrons. For neon, which has 10 electrons, we say that one electron sees an average potential due to the nucleus plus the other nine electrons. We imagine then that in the Schrödinger equation for each electron we put a $V(r)$ which is a $1/r$ field modified by a spherically symmetric charge density coming from the other electrons.

In this model each electron acts like an independent particle. The angular dependence of its wave function will be just the same as the ones we had for the hydrogen atom. There will be s -states, p -states, and so on; and they will have the various possible m -values. Since $V(r)$ no longer goes as $1/r$, the radial part of the wave functions will be somewhat different, but it will be qualitatively the same, so we will have the same radial quantum numbers, n . The energies of the states will also be somewhat different.

H

With these ideas, let's see what we get. The ground state of hydrogen has $l = m = 0$ and $n = 1$; we say the electron configuration is $1s$. The energy is -13.6 ev. This means that it takes 13.6 electron volts to pull the electron off the atom. We call this the "ionization energy", W_I . A large ionization energy means that it is harder to pull the electron off and, in general, that the material is chemically less active.

He

Now take helium. Both electrons can be in the same lowest state (one spin up and the other spin down). In this lowest state the electron moves in a potential which is for small r like a Coulomb field for $z = 2$ and for large r like a Coulomb field for $z = 1$. The result is a "hydrogen-like" $1s$ state with a somewhat lower energy. Both electrons occupy identical $1s$ states ($l = 0, m = 0$). The observed ionization energy (to remove *one* electron) is 24.6 electron volts. Since the $1s$ "shell" is now filled—we allow only two electrons—there is practically no tendency for an electron to be attracted from another atom. Helium is chemically inert.

Li

The lithium nucleus has a charge of 3. The electron states will again be hydrogen-like, and the three electrons will occupy the lowest three energy levels. Two will go into $1s$ states and the third will go into an $n = 2$ state. But with $l = 0$ or $l = 1$? In hydrogen these states have the same energy, but in other atoms they

don't, for the following reason. Remember that a $2s$ state has some amplitude to be near the nucleus while the $2p$ state does not. That means that a $2s$ electron will feel some of the triple electric charge of the Li nucleus, but that a $2p$ electron will stay out where the field looks like the Coulomb field of a single charge. The extra attraction lowers the energy of the $2s$ state relative to the $2p$ state. The energy levels will be roughly as shown in Fig. 19-8—which you should compare with the corresponding diagram for hydrogen in Fig. 19-7. So the lithium atom will have two electrons in $1s$ states and one in a $2s$. Since the $2s$ electron has a higher energy than a $1s$ electron it is relatively easily removed. The ionization energy of lithium is only 5.4 electron volts, and it is quite active chemically.

So you can see the patterns which develop; we have given in Table 19-2 a list of the first 36 elements, showing the states occupied by the electrons in the ground state of each atom. The Table gives the ionization energy for the most loosely bound electron, and the number of electrons occupying each “shell”—by which we mean states with the same n . Since the different l -states have different

Table 19-2

The electron configurations of the first 36 elements

Z	Element	$W_I(\text{ev})$	Electron Configuration									
			1s	2s	2p	3s	3p	3d	4s	4p	4d	4f
1	H <i>hydrogen</i>	13.6	1									
2	He <i>helium</i>	24.6	2									
3	Li <i>lithium</i>	5.4		1								
4	Be <i>beryllium</i>	9.3		2								
5	B <i>boron</i>	8.3		2	1							
6	C <i>carbon</i>	11.3	FILLED (2)	2	2	Number of electrons in each state						
7	N <i>nitrogen</i>	14.5		2	3							
8	O <i>oxygen</i>	13.6		2	4							
9	F <i>fluorine</i>	17.4		2	5							
10	Ne <i>neon</i>	21.6		2	6							
11	Na <i>sodium</i>	5.1				1						
12	Mg <i>magnesium</i>	7.6				2						
13	Al <i>aluminum</i>	6.0				2	1					
14	Si <i>silicon</i>	8.1	—FILLED—			2	2					
15	P <i>phosphorus</i>	10.5				2	3					
16	S <i>sulfur</i>	10.4	(2)	(8)		2	4					
17	Cl <i>chlorine</i>	13.0				2	5					
18	A <i>argon</i>	15.8				2	6					
19	K <i>potassium</i>	4.3							1			
20	Ca <i>calcium</i>	6.1							2			
21	Sc <i>scandium</i>	6.5						1	2			
22	Ti <i>titanium</i>	6.8						2	2			
23	V <i>vanadium</i>	6.7	—FILLED—					3	2			
24	Cr <i>chromium</i>	6.8						5	1			
25	Mn <i>manganese</i>	7.4	(2)	(8)		(8)		5	2			
26	Fe <i>iron</i>	7.9						6	2			
27	Co <i>cobalt</i>	7.9						7	2			
28	Ni <i>nickel</i>	7.6						8	2			
29	Cu <i>copper</i>	7.7						10	1			
30	Zn <i>zinc</i>	9.4						10	2			
31	Ga <i>gallium</i>	6.0							2	1		
32	Ge <i>germanium</i>	7.9	—FILLED—						2	2		
33	As <i>arsenic</i>	9.8							2	3		
34	Se <i>selenium</i>	9.7	(2)	(8)		(18)		2	4			
35	Br <i>bromine</i>	11.8						2	5			
36	Kr <i>krypton</i>	14.0						2	6			

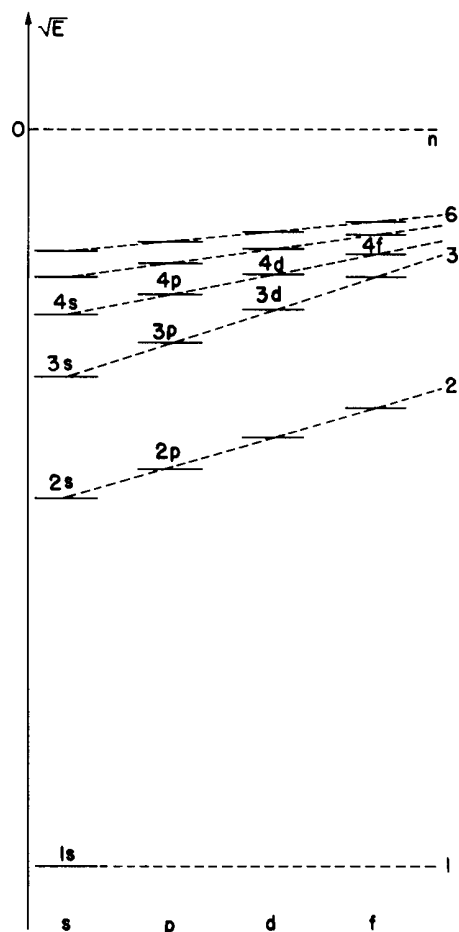


Fig. 19-8. Schematic energy level diagram for an atomic electron with other electrons present. (The scale is *not* the same as Fig. 19-7.)

energies, each l -value corresponds to a sub-shell of $2(2l + 1)$ possible states (of different m and electron spin). These all have the same energy—except for some very small effects we are neglecting.

Be

Beryllium is like lithium except that it has two electrons in the $2s$ state as well as two in the filled $1s$ shell.

B to Ne

Boron has 5 electrons. The fifth must go into a $2p$ state. There are $2 \times 3 = 6$ different $2p$ states, so we can keep adding electrons until we get to a total of 8. This takes us to neon. As we add these electrons we are also increasing Z , so the whole electron distribution gets pulled in closer and closer to the nucleus and the energy of the $2p$ states goes down. By the time we get to neon the ionization energy is up to 21.6 volts. Neon does not easily give up an electron. Also there are no more low-energy slots to be filled, so it won't try to grab an extra electron. Neon is chemically inert. Fluorine, on the other hand, *does* have an empty position where an electron can drop into a state of low energy, so it is quite active in chemical reactions.

Na to Ar

With sodium the eleventh electron must start a new shell—going into a $3s$ state. The energy level of this state is much higher; the ionization energy jumps down; and sodium is an active chemical. From sodium to argon the s and p states with $n = 3$ are occupied in exactly the same sequence as for lithium to neon. Angular configurations of the electrons in the outer unfilled shell have the same sequence, and the progression of ionization energies is quite similar. You can see why the chemical properties repeat with increasing atomic number. Magnesium acts chemically much like beryllium, silicon like carbon, and chlorine like fluorine. Argon is inert like neon.

You may have noticed that there is a slight peculiarity in the sequence of ionization energies between lithium and neon, and a similar one between sodium and argon. The last electron is bound to the oxygen atom somewhat less than we might expect. And sulphur is similar. Why should that be? We can understand it if we put in just a little bit of the effects of the interactions between individual electrons. Think of what happens when we put the first $2p$ electron onto the boron atom. It has six possibilities—three possible p -states, each with two spins. Imagine that the electron goes with spin up into the $m = 0$ state, which we have also called the “ z ” state because it hugs the z -axis. Now what will happen in carbon? There are now two $2p$ electrons. If one of them goes into the “ z ” state, where will the second one go? It will have lower energy if it stays away from the first electron, which it can do by going into, say, the “ x ” state of the $2p$ shell. (This state is, remember, just a linear combination of the $m = +1$ and $m = -1$ states.) Next, when we go to nitrogen, the three $2p$ electrons will have the lowest energy of mutual repulsion if they go one each into the “ x ,” “ y ,” and “ z ” configurations. For oxygen, however, the jig is up. The fourth electron must go into one of the filled states—with opposite spin. It is strongly repelled by the electron already in that state, so its energy will not be as low as it might otherwise be, and it is more easily removed. That explains the break in the sequence of binding energies which appears between nitrogen and oxygen, and between phosphorus and silicon.

K to Zn

After argon, you would, at first, think that the new electrons would start to fill up the $3d$ states. But they don't. As we described earlier—and illustrated in Fig. 19-7—the higher angular momentum states get pushed up in energy. By the time we get to the $3d$ states they are pushed to an energy a little bit above the energy of the $4s$ state. So in potassium the last electron goes into the $4s$ state. After this

shell is filled (with two electrons) at calcium, the $3d$ states begin to be filled for scandium, titanium, and vanadium.

The energies of the $3p$ and $4s$ states are so close together that small effects can shift the balance either way. By the time we get to put four electrons into the $3d$ states, their repulsion raises the energy of the $4s$ state just enough that its energy is slightly above the $3d$ energy, so one electron shifts over. For chromium we don't get a 4, 2 combination as we would have expected, but instead a 5, 1 combination. The new electron added to get manganese fills up the $4s$ shell again, and the states of the $3d$ shell are then occupied one by one until we reach copper.

Since the outermost shell of manganese, iron, cobalt, and nickel have the same configurations, however, they all tend to have similar chemical properties. (This effect is much more pronounced in the rare-earth elements which all have the same outer shell but a progressively filling inner shell which has much less influence on their chemical properties.)

In copper an electron is robbed from the $4s$ shell, finally completing the $3d$ shell. The energy of the 10, 1 combination is, however, so close to the 9, 2 configuration for copper that just the presence of another atom nearby can shift the balance. For this reason the two last electrons of copper are nearly equivalent, and copper can have a valence of either 1 or 2 (It sometimes acts as though its electrons were in the 9, 2 combination.) Similar things happen at other places and account for the fact that other metals, such as iron, combine chemically with either of two valences. By zinc, both the $3d$ and $4s$ shells are filled once and for all.

Ga to Kr

From gallium to krypton the sequence proceeds normally again, filling the $4p$ shell. The outer shells, the energies, and the chemical properties repeat the pattern of boron to neon and aluminum to argon.

Krypton, like argon and neon, is known as "noble" gas. All three are chemically "inert." This means only that, having filled shells of relatively low energy, there are few situations in which there is an energy advantage for them to join in a simple combination with other elements. Having a filled shell is not enough. Beryllium and magnesium have filled s -shells, but the energy of these shells is too high to lead to stability. Similarly, one would have expected another "noble" element at nickel, if the energy of the $3d$ shell had been lower (or the $4s$, higher). On the other hand, krypton is not completely inert; it will form a weakly-bound compound with chlorine.

Since our sample has turned up most of the main features of the periodic table, we stop our examination at element number 36—there are still seventy or so more!

We would like to bring up only one more point—that we not only can understand the valences to some extent but also can say something about the directional properties of the chemical bonds. Take an atom like oxygen which has four $2p$ electrons. The first three go into " x ," " y ," and " z " states and the fourth will double one of these states, leaving two—say " x " and " y "—vacant. Consider then what happens in H_2O . Each of the two hydrogens are willing to share an electron with the oxygen, helping the oxygen to fill a shell. These electrons will tend to go into the " x " and " y " vacancies. So the water molecule should have the two hydrogen atoms making a right angle with respect to the center of the oxygen. The angle is actually 105° . We can even understand why the angle is larger than 90° . In sharing their electrons the hydrogens end up with a net positive charge. The electric repulsion "strains" the wave functions and pushes the angle out to 105° . The same situation occurs in H_2S . But because the sulphur atom is larger, the two hydrogen atoms are farther apart, there is less repulsion, and the angle is only pushed out to about 93° . Selenium is even larger, so in H_2Se the angle is very nearly 90° .

We can use the same arguments to understand the geometry of ammonia, H_3N . Nitrogen has room for three more $2p$ electrons, one on each for the " x ," " y ," and " z " type states. The three hydrogens should join on at right angles to each other. The angles come out a little larger than 90° —again from the electric repul-

sion—but at least we see why the molecule of H_3N is not flat. The angles in phosphene, H_3P , are close to 90° , and in H_3As are still closer. We assumed that NH_3 was not flat when we described it as a two-state system. And the nonflatness is what makes the ammonia maser possible. Now we see that also that shape can be understood from our quantum mechanics.

The Schrödinger equation has been one of the great triumphs of physics. By providing the key to the underlying machinery of atomic structure it has given an explanation for atomic spectra, for chemistry, and for the nature of matter.

Operators

20-1 Operations and operators

All the things we have done so far in quantum mechanics could be handled with ordinary algebra, although we did from time to time show you some special ways of writing quantum-mechanical quantities and equations. We would like now to talk some more about some interesting and useful mathematical ways of describing quantum-mechanical things. There are many ways of approaching the subject of quantum mechanics, and most books use a different approach from the one we have taken. As you go on to read other books you might not see right away the connections of what you will find in them to what we have been doing. Although we will also be able to get a few useful results, the main purpose of this chapter is to tell you about some of the different ways of writing the same physics. Knowing them you should be able to understand better what other people are saying. When people were first working out classical mechanics they always wrote all the equations in terms of x -, y -, and z -components. Then someone came along and pointed out that all of the writing could be made much simpler by inventing the vector notation. It's true that when you come down to figuring something out you often have to convert the vectors back to their components. But it's generally much easier to see what's going on when you work with vectors and also easier to do many of the calculations. In quantum mechanics we were able to write many things in a simpler way by using the idea of the "state vector." The state vector $|\psi\rangle$ has, of course, nothing to do with geometric vectors in three dimensions but is an abstract symbol that *stands for a physical state*, identified by the "label," or "name," ψ . The idea is useful because the laws of quantum mechanics can be written as algebraic equations in terms of these symbols. For instance, our fundamental law that any state can be made up from a linear combination of base states is written as

$$|\psi\rangle = \sum_i C_i |i\rangle, \quad (20.1)$$

where the C_i are a set of ordinary (complex) numbers—the amplitudes $C_i = \langle i | \psi \rangle$ —while $|1\rangle$, $|2\rangle$, $|3\rangle$, and so on, stand for the base states in some base, or *representation*.

If you take some physical state and do something to it—like rotating it, or like waiting for the time Δt —you get a different state. We say, "performing an operation on a state produces a new state." We can express the same idea by an equation:

$$|\phi\rangle = \hat{A} |\psi\rangle. \quad (20.2)$$

An operation on a state produces another state. The *operator* \hat{A} stands for some particular operation. When this operation is performed on any state, say $|\psi\rangle$, it produces some other state $|\phi\rangle$.

What does Eq. (20.2) mean? We *define* it this way. If you multiply the equation by $\langle i |$ and expand $|\psi\rangle$ according to Eq. (20.1), you get

$$\langle i | \phi \rangle = \sum_j \langle i | \hat{A} | j \rangle \langle j | \psi \rangle. \quad (20.3)$$

(The states $|j\rangle$ are from the same set as $|i\rangle$.) This is now just an algebraic equation. The numbers $\langle i | \phi \rangle$ give the amount of each base state you will find in $|\phi\rangle$, and it is given in terms of a linear superposition of the amplitudes $\langle j | \psi \rangle$ that you find

20-1 Operations and operators

20-2 Average energies

20-3 The average energy of an atom

20-4 The position operator

20-5 The momentum operator

20-6 Angular momentum

20-7 The change of averages with time

$|\psi\rangle$ in each base state. The numbers $\langle i | \hat{A} | j \rangle$ are just the coefficients which tell how much of $\langle j | \psi \rangle$ goes into each sum. The operator \hat{A} is described numerically by the set of numbers, or “matrix,”

$$A_{ij} \equiv \langle i | \hat{A} | j \rangle. \quad (20.4)$$

So Eq. (20.2) is a high-class way of writing Eq. (20.3). Actually it is a little more than that; something more is implied. In Eq. (20.2) we do not make any reference to a set of base states. Equation (20.3) is an image of Eq. (20.2) in terms of some set of base states. But, as you know, you may use any set you wish. And this idea is implied in Eq. (20.2). The operator way of writing avoids making any particular choice. Of course, when you want to get definite you have to choose *some* set. When you make your choice, you use Eq. (20.3). So the *operator* equation (20.2) is a more abstract way of writing the *algebraic* equation (20.3). It's similar to the difference between writing

$$\mathbf{c} = \mathbf{a} \times \mathbf{b}$$

instead of

$$c_x = a_y b_z - a_z b_y,$$

$$c_y = a_z b_x - a_x b_z,$$

$$c_z = a_x b_y - a_y b_x.$$

The first way is much handier. When you want *results*, however, you will eventually have to give the components with respect to some set of axes. Similarly, if you want to be able to say what you really mean by \hat{A} , you will have to be ready to give the matrix A_{ij} in terms of *some* set of base states. So long as you have in mind some set A_{ij} , Eq. (20.2) means just the same as Eq. (20.3). (You should remember also that once you know a matrix for one particular set of base states you can always calculate the corresponding matrix that goes with any other base. You can transform the matrix from one “representation” to another.)

The operator equation in (20.2) also allows a new way of thinking. If we imagine some operator \hat{A} , we can use it with any state $|\psi\rangle$ to create a new state $\hat{A}|\psi\rangle$. Sometimes a “state” we get this way may be very peculiar—it may not represent any *physical* situation we are likely to encounter in nature. (For instance, we may get a state that is not normalized to represent one electron.) In other words, we may at times get “states” that are mathematically artificial. Such artificial “states” may still be useful, perhaps as the mid-point of some calculation.

We have already shown you many examples of quantum-mechanical operators. We have had the rotation operator $\hat{R}_y(\theta)$ which takes a state $|\psi\rangle$ and produces a new state, which is the old state as seen in a rotated coordinate system. We have had the parity (or inversion) operator \hat{P} , which makes a new state by reversing all coordinates. We have had the operators $\hat{\sigma}_x$, $\hat{\sigma}_y$, and $\hat{\sigma}_z$ for spin one-half particles.

The operator \hat{J}_z was defined in Chapter 17 in terms of the rotation operator for a small angle ϵ .

$$\hat{R}_z(\epsilon) = 1 + \frac{i}{\hbar} \epsilon \hat{J}_z. \quad (20.5)$$

This just means, of course, that

$$\hat{R}_z(\epsilon) |\psi\rangle = |\psi\rangle + \frac{i}{\hbar} \epsilon \hat{J}_z |\psi\rangle. \quad (20.6)$$

In this example, $\hat{J}_z |\psi\rangle$ is \hbar/ϵ times the state you get if you rotate $|\psi\rangle$ by the small angle ϵ and then subtract the original state. It represents a “state” which is the *difference* of two states.

One more example. We had an operator \hat{p}_x —called the momentum operator (x -component) defined in an equation like (20.6). If $\hat{D}_x(L)$ is the operator which

displaces a state along x by the distance L , then \hat{p}_x is defined by

$$\hat{D}_x(\delta) = 1 + \frac{i}{\hbar} \delta \hat{p}_x, \quad (20.7)$$

where δ is a small displacement. Displacing the state $|\psi\rangle$ along x by a small distance δ gives a new state $|\psi'\rangle$. We are saying that this new state is the old state plus a small new piece

$$\frac{i}{\hbar} \delta \hat{p}_x |\psi\rangle.$$

The operators we are talking about work on a state vector like $|\psi\rangle$, which is an abstract description of a physical situation. They are quite different from *algebraic* operators which work on mathematical functions. For instance, d/dx is an “operator” that works on $f(x)$ by changing it to a new function $f'(x) = df/dx$. Another example is the algebraic operator ∇^2 . You can see why the same word is used in both cases, but you should keep in mind that the two kinds of operators are different. A quantum-mechanical operator \hat{A} does *not* work on an algebraic function, but on a state vector like $|\psi\rangle$. Both kinds of operators are used in quantum mechanics and often in similar kinds of equations, as you will see a little later. When you are first learning the subject it is well to keep the distinction always in mind. Later on, when you are more familiar with the subject, you will find that it is less important to keep any sharp distinction between the two kinds of operators. You will, indeed, find that most books generally use the same notation for both!

We’ll go on now and look at some useful things you can do with operators. But first, one special remark. Suppose we have an operator \hat{A} whose matrix in some base is $A_{ij} \equiv \langle i | \hat{A} | j \rangle$. The amplitude that the state $\hat{A} |\psi\rangle$ is also in some other state $|\phi\rangle$ is $\langle \phi | \hat{A} | \psi \rangle$. Is there some meaning to the complex conjugate of this amplitude? You should be able to show that

$$\langle \phi | \hat{A} | \psi \rangle^* = \langle \psi | \hat{A}^\dagger | \phi \rangle, \quad (20.8)$$

where \hat{A}^\dagger (read “A dagger”) is an operator whose matrix elements are

$$A_{ij}^\dagger = (A_{ji})^*. \quad (20.9)$$

To get the i, j element of A^\dagger you go to the j, i element of \hat{A} (the indexes are reversed) and take its complex conjugate. The amplitude that the state $\hat{A}^\dagger |\phi\rangle$ is in $|\psi\rangle$ is the complex conjugate of the amplitude that $\hat{A} |\psi\rangle$ is in $|\phi\rangle$. The operator \hat{A}^\dagger is called the “Hermitian adjoint” of \hat{A} . Many important operators of quantum mechanics have the special property that when you take the Hermitian adjoint, you get the same operator back. If \hat{B} is such an operator, then

$$\hat{B}^\dagger = \hat{B},$$

and it is called a “self-adjoint” or “Hermitian,” operator.

20-2 Average energies

So far we have reminded you mainly of what you already know. Now we would like to discuss a new question. How would you find the *average* energy of a system—say, an atom? If an atom is in a particular state of definite energy and you measure the energy, you will find a certain energy E . If you keep repeating the measurement on each one of a whole series of atoms which are all selected to be in the same state, all the measurements will give E , and the “average” of your measurements will, of course, be just E .

Now, however, what happens if you make the measurement on some state $|\psi\rangle$ which is *not* a stationary state? Since the system does not have a definite energy, one measurement would give one energy, the same measurement on another atom in the same state would give a different energy, and so on. What would you get for the average of a whole series of energy measurements?

We can answer the question by projecting the state $|\psi\rangle$ onto the set of states of definite energy. To remind you that this is a special base set, we'll call the states $|\eta_i\rangle$. Each of the states $|\eta_i\rangle$ has a definite energy E_i . In this representation,

$$|\psi\rangle = \sum_i C_i |\eta_i\rangle. \quad (20.10)$$

When you make an energy measurement and get some number E_i , you have found that the system was in the state η_i . But you may get a different number for each measurement. Sometimes you will get E_1 , sometimes E_2 , sometimes E_3 , and so on. The *probability* that you observe the energy E_1 is just the probability of finding the system in the state $|\eta_1\rangle$, which is, of course, just the absolute square of the amplitude $C_1 = \langle\eta_1|\psi\rangle$. The probability of finding each of the possible energies E_i is

$$P_i = |C_i|^2. \quad (20.11)$$

How are these probabilities related to the mean value of a whole sequence of energy measurements? Let's imagine that we get a series of measurements like this: $E_1, E_7, E_{11}, E_9, E_1, E_{10}, E_7, E_2, E_3, E_9, E_6, E_4$, and so on. We continue for, say, a thousand measurements. When we are finished we add all the energies and divide by one thousand. That's what we mean by the average. There's also a short-cut to adding all the numbers. You can count up how many times you get E_1 , say that is N_1 , and then count up the number of times you get E_2 , call that N_2 , and so on. The sum of all the energies is certainly just

$$N_1 E_1 + N_2 E_2 + N_3 E_3 + \cdots = \sum_i N_i E_i.$$

The average energy is this sum divided by the total number of measurements which is just the sum of all the N_i 's, which we can call N ;

$$E_{\text{av}} = \frac{\sum_i N_i E_i}{N}. \quad (20.12)$$

We are almost there. What we *mean* by the probability of something happening is just the number of times we expect it to happen divided by the total number of tries. The ratio N_i/N should—for large N —be very near to P_i , the probability of finding the state $|\eta_i\rangle$, although it will not be exactly P_i because of the statistical fluctuations. Let's write the predicted (or “expected”) average energy as $\langle E \rangle_{\text{av}}$; then we can say that

$$\langle E \rangle_{\text{av}} = \sum_i P_i E_i. \quad (20.13)$$

The same arguments apply for any measurement. The average value of a measured quantity A should be equal to

$$\langle A \rangle_{\text{av}} = \sum_i P_i A_i,$$

where A_i are the various possible values of the observed quantity, and P_i is the probability of getting that value.

Let's go back to our quantum-mechanical state $|\psi\rangle$. Its average energy is

$$\langle E \rangle_{\text{av}} = \sum_i |C_i|^2 E_i = \sum_i C_i^* C_i E_i. \quad (20.14)$$

Now watch this trickery! First, we write the sum as

$$\sum_i \langle\psi|\eta_i\rangle E_i \langle\eta_i|\psi\rangle. \quad (20.15)$$

Next we treat the left-hand $\langle\psi|$ as a common “factor.” We can take this factor out of the sum, and write it as

$$\langle\psi| \left\{ \sum_i |\eta_i\rangle E_i \langle\eta_i| \right\} |\psi\rangle.$$

This expression has the form

$$\langle \psi | \phi \rangle,$$

where $|\phi\rangle$ is some “cooked-up” state defined by

$$|\phi\rangle = \sum_i |\eta_i\rangle E_i \langle \eta_i | \psi \rangle. \quad (20.16)$$

It is, in other words, the state you get if you take each base state $|\eta_i\rangle$ in the amount $E_i \langle \eta_i | \psi \rangle$.

Now remember what we mean by the states $|\eta_i\rangle$. They are supposed to be the stationary states—by which we mean that for each one,

$$\hat{H} |\eta_i\rangle = E_i |\eta_i\rangle.$$

Since E_i is just a number, the right-hand side is the same as $|\eta_i\rangle E_i$, and the sum in Eq. (20.16) is the same as

$$\sum_i \hat{H} |\eta_i\rangle \langle \eta_i | \psi \rangle.$$

Now i appears only in the famous combination that contracts to unity, so

$$\sum_i \hat{H} |\eta_i\rangle \langle \eta_i | \psi \rangle = \hat{H} \sum_i |\eta_i\rangle \langle \eta_i | \psi \rangle = \hat{H} |\psi\rangle.$$

Magic! Equation (20.16) is the same as

$$|\phi\rangle = \hat{H} |\psi\rangle. \quad (20.17)$$

The average energy of the state $|\psi\rangle$ can be written very prettily as

$$\langle E \rangle_{\text{av}} = \langle \psi | \hat{H} | \psi \rangle. \quad (20.18)$$

To get the average energy you operate on $|\psi\rangle$ with \hat{H} , and then multiply by $\langle \psi |$. A simple result.

Our new formula for the average energy is not only pretty. It is also useful, because now we don't need to say anything about any particular set of base states. We don't even have to know all of the possible energy levels. When we go to calculate, we'll need to describe our state in terms of *some* set of base states, but if we know the Hamiltonian matrix H_{ij} for *that* set we can get the average energy. Equation (19.18) says that for *any* set of base states $|i\rangle$, the average energy can be calculated from

$$\langle E \rangle_{\text{av}} = \sum_{i,j} \langle \psi | i \rangle \langle i | \hat{H} | j \rangle \langle j | \psi \rangle, \quad (20.19)$$

where the amplitudes $\langle i | \hat{H} | j \rangle$ are just the elements of the matrix H_{ij} .

Let's check this result for the special case that the states $|i\rangle$ are the definite energy states. For them, $\hat{H} |j\rangle = E_j |j\rangle$, so $\langle i | \hat{H} | j \rangle = E_j \delta_{ij}$ and

$$\langle E \rangle_{\text{av}} = \sum_{i,j} \langle \psi | i \rangle E_i \delta_{ij} \langle j | \psi \rangle = \sum_i E_i \langle \psi | i \rangle \langle i | \psi \rangle,$$

which is right.

Equation (20.19) can, incidentally, be extended to other physical measurements which you can express as an operator. For instance, \hat{L}_z is the operator of the z-component of the angular momentum L . The average of the z-component for the state $|\psi\rangle$ is

$$\langle L_z \rangle_{\text{av}} = \langle \psi | \hat{L}_z | \psi \rangle.$$

One way to prove it is to think of some situation in which the energy is proportional to the angular momentum. Then all the arguments go through in the same way.

In summary, if a physical observable A is related to a suitable quantum-mechanical operator \hat{A} , the average value of A for the state $|\psi\rangle$ is given by

$$\langle A \rangle_{\text{av}} = \langle \psi | \hat{A} | \psi \rangle. \quad (20.20)$$

By this we mean that

$$A_{\text{av}} = \langle \psi | \phi \rangle, \quad (20.21)$$

with

$$|\phi\rangle = \hat{A} |\psi\rangle. \quad (20.22)$$

20-3 The average energy of an atom

Suppose we want the average energy of an atom in a state described by a wave function $\psi(\mathbf{r})$; How do we find it? Let's first think of a one-dimensional situation with a state $|\psi\rangle$ defined by the amplitude $\langle x | \psi \rangle = \psi(x)$. We are asking for the special case of Eq. (20.19) applied to the coordinate representation. Following our usual procedure, we replace the states $|i\rangle$ and $|j\rangle$ by $|x\rangle$ and $|x'\rangle$, and change the sums to integrals. We get

$$\langle E \rangle_{\text{av}} = \int \int \langle \psi | x \rangle \langle x | \hat{H} | x' \rangle \langle x' | \psi \rangle dx dx'. \quad (20.23)$$

This integral can, if we wish, be written in the following way:

$$\int \langle \psi | x \rangle \langle x | \phi \rangle dx, \quad (20.24)$$

with

$$\langle x | \phi \rangle = \int \langle x | \hat{H} | x' \rangle \langle x' | \psi \rangle dx'. \quad (20.25)$$

The integral over x' in (20.25) is the same one we had in Chapter 16—see Eq. (16.50) and Eq. (16.52)—and is equal to

$$-\frac{\hbar^2}{2m} \frac{d^2}{dx^2} \psi(x) + V(x)\psi(x).$$

We can therefore write

$$\langle x | \phi \rangle = \left\{ -\frac{\hbar^2}{2m} \frac{d^2}{dx^2} + V(x) \right\} \psi(x). \quad (20.26)$$

Remember that $\langle \psi | x \rangle = \langle x | \psi \rangle^* = \psi^*(x)$; using this equality, the average energy in Eq. (20.23) can be written as

$$\langle E \rangle_{\text{av}} = \int \psi^*(x) \left\{ -\frac{\hbar^2}{2m} \frac{d^2}{dx^2} + V \right\} \psi(x) dx. \quad (20.27)$$

Given a wave function $\psi(x)$, you can get the average energy by doing this integral. You can begin to see how we can go back and forth from the state-vector ideas to the wave-function ideas.

The quantity in the braces of Eq. (20.27) is an *algebraic* operator.† We will write it as $\hat{\mathcal{H}}$

$$\hat{\mathcal{H}} = -\frac{\hbar^2}{2m} \frac{d^2}{dx^2} + V.$$

With this notation Eq. (20.23) becomes

$$\langle E \rangle_{\text{av}} = \int \psi^*(x) \hat{\mathcal{H}} \psi(x) dx. \quad (20.28)$$

The algebraic operator $\hat{\mathcal{H}}$ defined here is, of course, not identical to the quantum-mechanical operator \hat{H} . The new operator works on a function of position $\psi(x) = \langle x | \psi \rangle$ to give a new function of x , $\phi(x) = \langle x | \phi \rangle$; while \hat{H}

† The “operator” $V(x)$ means “multiply by $V(x)$.”

operates on a state vector $|\psi\rangle$ to give another state vector $|\phi\rangle$, without implying the coordinate representation or any particular representation at all. Nor is $\hat{\mathcal{C}}$ strictly the same as \hat{H} even in the coordinate representation. If we choose to work in the coordinate representation, we would interpret \hat{H} in terms of a matrix $\langle x | \hat{H} | x' \rangle$ which depends somehow on the two "indices" x and x' ; that is, we expect—according to Eq. (20.25)—that $\langle x | \phi \rangle$ is related to all the amplitudes $\langle x | \psi \rangle$ by an integration. On the other hand, we find that $\hat{\mathcal{C}}$ is a differential operator. We have already worked out in Section 16-5 the connection between $\langle x | \hat{H} | x' \rangle$ and the algebraic operator $\hat{\mathcal{C}}$.

We should make one qualification on our results. We have been assuming that the amplitude $\psi(x) = \langle x | \psi \rangle$ is normalized. By this we mean that the scale has been chosen so that

$$\int |\psi(x)|^2 dx = 1;$$

so the probability of finding the electron *somewhere* is unity. If you should choose to work with a $\psi(x)$ which is not normalized you should write

$$\langle E \rangle_{\text{av}} = \frac{\int \psi^*(x) \hat{\mathcal{C}} \psi(x) dx}{\int \psi^*(x) \psi(x) dx}. \quad (20.29)$$

It's the same thing.

Notice the similarity in form between Eq. (20.28) and Eq. (20.18). These two ways of writing the same result appear often when you work with the x -representation. You can go from the first form to the second with any \hat{A} which is a *local* operator, where a local operator is one which in the integral

$$\int \langle x | \hat{A} | x' \rangle \langle x' | \psi \rangle dx'$$

can be written as $\hat{\alpha} \psi(x)$, where $\hat{\alpha}$ is a differential algebraic operator. There are, however, operators for which this is not true. For them you must work with the basic equations in (20.21) and (20.22).

You can easily extend the derivation to three dimensions. The result is that†

$$\langle E \rangle_{\text{av}} = \int \psi(\mathbf{r}) \hat{\mathcal{C}} \psi(\mathbf{r}) d \text{Vol}, \quad (20.30)$$

with

$$\hat{\mathcal{C}} = -\frac{\hbar^2}{2m} \nabla^2 + V(\mathbf{r}), \quad (20.31)$$

and with the understanding that

$$\int |\psi|^2 d \text{Vol} = 1. \quad (20.32)$$

The same equations can be extended to systems with several electrons in a fairly obvious way, but we won't bother to write down the results.

With Eq. (20.30) we can calculate the average energy of an atomic state even without knowing its energy levels. All we need is the wave function. It's an important law. We'll tell you about one interesting application. Suppose you want to know the ground-state energy of some system—say the helium atom, but it's too hard to solve Schrödinger's equation for the wave function, because there are too many variables. Suppose, however, that you take a guess at the wave function—pick any function you like—and calculate the average energy. That is, you use Eq. (20.29)—generalized to three dimensions—to find what the average energy would be if the atom were really in the state described by this wave function. This energy will certainly be higher than the ground-state energy which is the lowest

† We write $d \text{Vol}$ for the element of volume. It is, of course, just $dx dy dz$, and the integral goes from $-\infty$ to $+\infty$ in all three coordinates.

possible energy the atom can have.† Now pick another function and calculate its average energy. If it is lower than your first choice you are getting closer to the true ground-state energy. If you keep on trying all sorts of artificial states you will be able to get lower and lower energies, which come closer and closer to the ground-state energy. If you are clever, you will try some functions which have a few adjustable parameters. When you calculate the energy it will be expressed in terms of these parameters. By varying the parameters to give the lowest possible energy, you are trying out a whole class of functions at once. Eventually you will find that it is harder and harder to get lower energies and you will begin to be convinced that you are fairly close to the lowest possible energy. The helium atom has been solved in just this way—not by solving a differential equation, but by making up a special function with a lot of adjustable parameters which are eventually chosen to give the lowest possible value for the average energy.

20-4 The position operator

What is the average value of the position of an electron in an atom? For any particular state $|\psi\rangle$ what is the average value of the coordinate x ? We'll work in one dimension and let you extend the ideas to three dimensions or to systems with more than one particle. We have a state described by $\psi(x)$, and we keep measuring x over and over again. What is the average? It is

$$\int xP(x) dx,$$

where $P(x)$ is the probability of finding the electron in a little element dx at x . Suppose the probability density $P(x)$ varies with x as shown in Fig. 20-1. The electron is most likely to be found near the peak of the curve. The average value of x is also somewhere near the peak. It is, in fact, just the center of gravity of the area under the curve.

We have seen earlier that $P(x)$ is just $|\psi(x)|^2 = \psi^*(x)\psi(x)$, so we can write the average of x as

$$\langle x \rangle_{av} = \int \psi^*(x)x\psi(x) dx. \quad (20.33)$$

Our equation for $\langle x \rangle_{av}$ has the same form as Eq. (20.33). For the average energy, the energy operator $\hat{\mathcal{E}}$ appears between the two ψ 's, for the average position there is just x . (If you wish you can consider x to be the algebraic operator "multiply by x ."*) We can carry the parallelism still further, expressing the average position in a form which corresponds to Eq. (20.18). Suppose we just write

$$\langle x \rangle_{av} = \langle \psi | \alpha \rangle \quad (20.34)$$

with

$$| \alpha \rangle = \hat{x} | \psi \rangle, \quad (20.35)$$

and then see if we can find the operator \hat{x} which generates the state $|\alpha\rangle$, which will make Eq. (20.34) agree with Eq. (20.33). That is, we must find a $|\alpha\rangle$, so that

$$\langle \psi | \alpha \rangle = \langle x \rangle_{av} = \int \langle \psi | x \rangle x \langle x | \psi \rangle dx. \quad (20.36)$$

First, let's expand $\langle \psi | \phi \rangle$ in the x -representation. It is

$$\langle \psi | \alpha \rangle = \int \langle \psi | x \rangle \langle x | \alpha \rangle dx. \quad (20.37)$$

Now compare the integrals in the last two equations. You see that in the x -representation

$$\langle x | \alpha \rangle = x \langle x | \psi \rangle. \quad (20.38)$$

† You can also look at it this way. Any function (that is, state) you choose can be written as a linear combination of the base states which are definite energy states. Since in this combination there is a mixture of higher energy states in with the lowest energy state, the average energy will be higher than the ground-state energy.

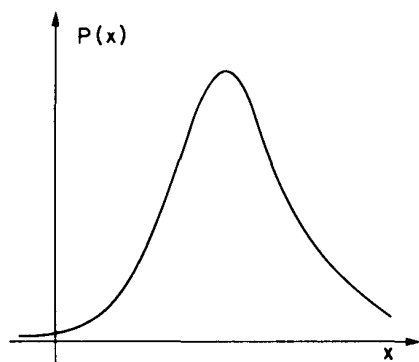


Fig. 20-1. A curve of probability density representing a localized particle.

Operating on $|\psi\rangle$ with \hat{x} to get $|\alpha\rangle$ is equivalent to multiplying $\psi(x) = \langle x|\psi\rangle$ by x to get $\alpha(x) = \langle x|\alpha\rangle$. We have a definition of \hat{x} in the coordinate representation.[†]

[We have not bothered to try to get the x -representation of the matrix of the operator \hat{x} . If you are ambitious you can try to show that

$$\langle x|\hat{x}|x'\rangle = x\delta(x - x'). \quad (20.39)$$

You can then work out the amusing result that

$$\hat{x}|x\rangle = x|x\rangle. \quad (20.40)$$

The operator \hat{x} has the interesting property that when it works on the base states $|x\rangle$ it is equivalent to multiplying by x .]

Do you want to know the average value of x^2 ? It is

$$\langle x^2 \rangle_{\text{av}} = \int \psi^*(x)x^2\psi(x) dx. \quad (20.41)$$

Or, if you prefer you can write

$$\langle x^2 \rangle_{\text{av}} = \langle \psi|\alpha'\rangle$$

with

$$|\alpha'\rangle = \hat{x}^2|\psi\rangle. \quad (20.42)$$

By \hat{x}^2 we mean $\hat{x}\hat{x}$ —the two operators are used one after the other. With the second form you can calculate $\langle x^2 \rangle_{\text{av}}$ using any representation (base-states) you wish. If you want the average of x^n , or of any polynomial in x , you can see how to get it.

20-5 The momentum operator

Now we would like to calculate the mean *momentum* of an electron—again, we'll stick to one dimension. Let $P(p)dp$ be the probability that a measurement will give a momentum between p and $p + dp$. Then

$$\langle p \rangle_{\text{av}} = \int p P(p) dp. \quad (20.43)$$

Now we let $\langle p|\psi\rangle$ be the amplitude that the state $|\psi\rangle$ is in a definite momentum state $|p\rangle$. This is the same amplitude we called $\langle \text{mom } p|\psi\rangle$ in Section 16-3 and is a function of p just as $\langle x|\psi\rangle$ is a function of x . There we chose to normalize the amplitude so that

$$P(p) = \frac{1}{2\pi\hbar} |\langle p|\psi\rangle|^2. \quad (20.44)$$

We have, then,

$$\langle p \rangle_{\text{av}} = \int \langle \psi|p\rangle p \langle p|\psi\rangle \frac{dp}{2\pi\hbar}. \quad (20.45)$$

The form is quite similar to what we had for $\langle x \rangle_{\text{av}}$.

If we want, we can play exactly the same game we did with $\langle x \rangle_{\text{av}}$. First, we can write the integral above as

$$\int \langle \psi|p\rangle \langle p|\beta\rangle \frac{dp}{2\pi\hbar}. \quad (20.46)$$

You should now recognize this equation as just the expanded form of the amplitude $\langle \psi|\beta\rangle$ —expanded in terms of the base states of definite momentum. From Eq.

[†] Equation (20.38) does *not* mean that $|\alpha\rangle = x|\psi\rangle$. You cannot “factor out” the $\langle x|$, because the multiplier x in front of $\langle x|\psi\rangle$ is a number which is different for each state $\langle x|$. It is the value of the coordinate of the electron in the state $|x\rangle$. See Eq. (20.40).

(20.45) the state $|\beta\rangle$ is defined *in the momentum representation* by

$$\langle p | \beta \rangle = p \langle p | \psi \rangle \quad (20.47)$$

That is, we can now write

$$\langle p \rangle_{\text{av}} = \langle \psi | \beta \rangle \quad (20.48)$$

with

$$|\beta\rangle = \hat{p} |\psi\rangle, \quad (20.49)$$

where the operator \hat{p} is defined in terms of the p -representation by Eq. (20.47).

[Again, you can if you wish show that the matrix form of \hat{p} is

$$\langle p | \hat{p} | p' \rangle = p \delta(p - p'), \quad (20.50)$$

and that

$$\hat{p} | p \rangle = p | p \rangle. \quad (20.51)$$

It works out the same as for x .]

Now comes an interesting question. We can write $\langle p \rangle_{\text{av}}$ as we have done in Eqs. (20.45) and (20.48), and we know the meaning of the operator \hat{p} *in the momentum representation*. But how should we interpret \hat{p} in the *coordinate* representation? That is what we will need to know if we have some wave function $\psi(x)$, and we want to compute its average momentum. Let's make clear what we mean. If we start by saying that $\langle p \rangle_{\text{av}}$ is given by Eq. (20.48), we can expand that equation in terms of the p -representation to get back to Eq. (20.45). If we are given the p -description of the state—namely the amplitude $\langle p | \psi \rangle$, which is an algebraic function of the momentum p —we can get $\langle p | \phi \rangle$ from Eq. (20.47) and proceed to evaluate the integral. The question now is: What do we do if we are given a description of the state in the x -representation, namely the wave function $\psi(x) = \langle x | \psi \rangle$?

Well, let's start by expanding Eq. (20.48) in the x -representation. It is

$$\langle p \rangle_{\text{av}} = \int \langle \psi | x \rangle \langle x | \beta \rangle dx. \quad (20.52)$$

Now, however, we need to know what the state $|\beta\rangle$ is in the x -representation. If we can find it, we can carry out the integral. So our problem is to find the function $\beta(x) = \langle x | \beta \rangle$.

We can find it in the following way. In Section 16-3 we saw how $\langle p | \beta \rangle$ was related to $\langle x | \beta \rangle$. According to Eq. (16.24),

$$\langle p | \beta \rangle = \int e^{-ipx/\hbar} \langle x | \beta \rangle dx. \quad (20.53)$$

If we know $\langle p | \beta \rangle$ we can solve this equation for $\langle x | \beta \rangle$. What we want, of course, is to express the result somehow in terms of $\psi(x) = \langle x | \psi \rangle$, which we are assuming to be known. Suppose we start with Eq. (20.47) and again use Eq. (16.24) to write

$$\langle p | \beta \rangle = p \langle p | \psi \rangle = p \int e^{-ipx/\hbar} \psi(x) dx. \quad (20.54)$$

Since the integral is over x we can put the p inside the integral and write

$$\langle p | \beta \rangle = \int e^{-ipx/\hbar} p \psi(x) dx. \quad (20.55)$$

Compare this with (20.53). You would say that $\langle x | \beta \rangle$ is equal to $p\psi(x)$. No, No! The wave function $\langle x | \beta \rangle = \beta(x)$ can depend only on x —not on p . That's the whole problem.

However, some ingenious fellow discovered that the integral in (20.55) could be integrated by parts. The derivative of $e^{-ipx/\hbar}$ with respect to x is $(-i/\hbar)pe^{-ipx/\hbar}$, so the integral in (20.55) is equivalent to

$$-\frac{\hbar}{i} \int \frac{d}{dx} (e^{-ipx/\hbar}) \psi(x) dx.$$

If we integrate by parts, it becomes

$$-\frac{\hbar}{i} [e^{-ipx/\hbar} \psi(x)]_{-\infty}^{+\infty} + \frac{\hbar}{i} \int e^{-ipx/\hbar} \frac{d\psi}{dx} dx.$$

So long as we are considering bound states, so that $\psi(x)$ goes to zero at $x = \pm \infty$, the bracket is zero and we have

$$\langle p | \beta \rangle = \frac{\hbar}{i} \int e^{-ipx/\hbar} \frac{d\psi}{dx} dx. \quad (20.56)$$

Now compare this result with Eq. (20.53). You see that

$$\langle x | \beta \rangle = \frac{\hbar}{i} \frac{d}{dx} \psi(x). \quad (20.57)$$

We have the necessary piece to be able to complete Eq. (20.52). The answer is

$$\langle p \rangle_{av} = \int \psi^*(x) \frac{\hbar}{i} \frac{d}{dx} \psi(x) dx. \quad (20.58)$$

We have found how Eq. (20.48) looks in the coordinate representation.

Now you should begin to see an interesting pattern developing. When we asked for the average energy of the state $|\psi\rangle$ we said it was

$$\langle E \rangle_{av} = \langle \psi | \phi \rangle, \quad \text{with} \quad |\phi\rangle = \hat{H} |\psi\rangle.$$

The same thing is written in the coordinate world as

$$\langle E \rangle_{av} = \int \psi^*(x) \phi(x) dx \quad \text{with} \quad \phi(x) = \hat{\mathcal{H}} \psi(x).$$

Here $\hat{\mathcal{H}}$ is an *algebraic* operator which works a function of x . When we asked about the average value of x , we found that it could also be written

$$\langle x \rangle_{av} = \langle \psi | \alpha \rangle, \quad \text{with} \quad |\alpha\rangle = \hat{x} |\psi\rangle.$$

In the coordinate world the corresponding equations are

$$\langle x \rangle_{av} = \int \psi^*(x) \alpha(x) dx, \quad \text{with} \quad \alpha(x) = x\psi(x).$$

When we asked about the average value of p , we wrote

$$\langle p \rangle_{av} = \langle \psi | \beta \rangle, \quad \text{with} \quad |\beta\rangle = \hat{p} |\psi\rangle.$$

In the coordinate world the equivalent equations were

$$\langle p \rangle_{av} = \int \psi(x) \beta(x) dx, \quad \text{with} \quad \beta(x) = \frac{\hbar}{i} \frac{d}{dx} \psi(x).$$

In each of our three examples we start with the state $|\psi\rangle$ and produce another (hypothetical) state by a *quantum-mechanical* operator. In the coordinate representation we generate the corresponding wave function by operating on the wave function $\psi(x)$ with an *algebraic* operator. There are the following one-to-one correspondences (for one-dimensional problems):

$$\begin{aligned} \hat{H} &\rightarrow \hat{\mathcal{H}} = -\frac{\hbar^2}{2m} \frac{d^2}{dx^2} + V(x), \\ \hat{x} &\rightarrow x, \\ \hat{p}_x &\rightarrow \hat{\mathcal{P}}_x = \frac{\hbar}{i} \frac{\partial}{\partial x}. \end{aligned} \quad (20.59)$$

Table 20-1

Physical Quantity	Operator	Coordinate Form
Energy	\hat{H}	$\hat{\mathcal{H}} = -\frac{\hbar^2}{2m} \nabla^2 + V(\mathbf{r})$
Position	\hat{x}	x
	\hat{y}	y
	\hat{z}	z
Momentum	\hat{p}_x	$\hat{\mathcal{P}}_x = \frac{\hbar}{i} \frac{\partial}{\partial x}$
	\hat{p}_y	$\hat{\mathcal{P}}_y = \frac{\hbar}{i} \frac{\partial}{\partial y}$
	\hat{p}_z	$\hat{\mathcal{P}}_z = \frac{\hbar}{i} \frac{\partial}{\partial z}$

In this list, we have introduced the symbol $\hat{\mathcal{P}}_x$ for the algebraic operator $(\hbar/i)\partial/\partial x$:

$$\hat{\mathcal{P}}_x = \frac{\hbar}{i} \frac{\partial}{\partial x}, \quad (20.60)$$

and we have inserted the x subscript on $\hat{\mathcal{P}}$ to remind you that we have been working only with the x -component of momentum.

You can easily extend the results to three dimensions. For the other components of the momentum,

$$\hat{p}_y \rightarrow \hat{\mathcal{P}}_y = \frac{\hbar}{i} \frac{\partial}{\partial y},$$

$$\hat{p}_z \rightarrow \hat{\mathcal{P}}_z = \frac{\hbar}{i} \frac{\partial}{\partial z}.$$

If you want, you can even think of an operator of the *vector* momentum and write

$$\hat{\mathbf{p}} \rightarrow \hat{\mathcal{P}} = \frac{\hbar}{i} \left(\mathbf{e}_x \frac{\partial}{\partial x} + \mathbf{e}_y \frac{\partial}{\partial y} + \mathbf{e}_z \frac{\partial}{\partial z} \right),$$

where \mathbf{e}_x , \mathbf{e}_y , and \mathbf{e}_z are the unit vectors in the three directions. It looks even more elegant if we write

$$\hat{\mathbf{p}} \rightarrow \hat{\mathcal{P}} = \frac{\hbar}{i} \nabla. \quad (20.61)$$

Our general result is that for at least some quantum-mechanical operators, there are corresponding algebraic operators in the coordinate representation. We summarize our results so far—extended to three dimensions—in Table 20-1. For each operator we have the two equivalent forms:‡

$$|\phi\rangle = \hat{A} |\psi\rangle \quad (20.62)$$

or

$$\phi(\mathbf{r}) = \hat{A}\psi(\mathbf{r}). \quad (20.63)$$

We will now give a few illustrations of the use of these ideas. The first one is just to point out the relation between $\hat{\mathcal{P}}$ and $\hat{\mathcal{H}}$. If we use $\hat{\mathcal{P}}_x$ twice, we get

$$\hat{\mathcal{P}}_x \hat{\mathcal{P}}_x = -\hbar^2 \frac{\partial^2}{\partial x^2}.$$

‡ In many books the same symbol is used for \hat{A} and $\hat{\mathcal{A}}$, because they both stand for the same physics, and because it is convenient not to have to write different kinds of letters. You can usually tell which one is intended by the context.

This means that we can write the equality

$$\hat{\mathcal{H}} = \frac{1}{2m} \{ \hat{p}_x \hat{p}_x + \hat{p}_y \hat{p}_y + \hat{p}_z \hat{p}_z \} + V(\mathbf{r}).$$

Or, using the vector notation,

$$\hat{\mathcal{H}} = \frac{1}{2m} \hat{\mathbf{p}} \cdot \hat{\mathbf{p}} + V(\mathbf{r}). \quad (20.64)$$

(In an algebraic operator, any term without the operator symbol ($\hat{}$) means just a straight multiplication.) This equation is nice because it's easy to remember if you haven't forgotten your classical physics. Everyone knows that the energy is (nonrelativistically) just the kinetic energy $p^2/2m$ plus the potential energy, and $\hat{\mathcal{H}}$ is the operator of the total energy.

This result has impressed people so much that they try to teach students all about classical physics before quantum mechanics. (We think differently!) But such parallels are often misleading. For one thing, when you have operators, the *order* of various factors is important; but that is not true for the factors in a classical equation.

In Chapter 17 we defined an operator \hat{p}_x in terms of the displacement operator \hat{D}_x by [see Eq. (17.27)]

$$|\psi'\rangle = \hat{D}_x(\delta) |\psi\rangle = \left(1 + \frac{i}{\hbar} \hat{p}_x \delta \right) |\psi\rangle, \quad (20.65)$$

where δ is a *small* displacement. We should show you that this is equivalent to our new definition. According to what we have just worked out, this equation should mean the same as

$$\psi'(x) = \psi(x) + \frac{\partial \psi}{\partial x} \delta$$

But the right-hand side is just the Taylor expansion of $\psi(x + \delta)$, which is certainly what you get if you displace the state to the left by δ (or shift the coordinates to the right by the same amount). Our two definitions of \hat{p} agree!

Let's use this fact to show something else. Suppose we have a bunch of particles which we label 1, 2, 3, . . . , in some complicated system. (To keep things simple we'll stick to one dimension.) The wave function describing the state is a function of all the coordinates x_1, x_2, x_3, \dots . We can write it as $\psi(x_1, x_2, x_3, \dots)$. Now displace the system (to the left) by δ . The new wave function

$$\psi'(x_1, x_2, x_3, \dots) = \psi(x_1 + \delta, x_2 + \delta, x_3 + \delta, \dots)$$

can be written as

$$\begin{aligned} \psi'(x_1, x_2, x_3, \dots) &= \psi(x_1, x_2, x_3, \dots) \\ &+ \left\{ \delta \frac{\partial \psi}{\partial x_1} + \delta \frac{\partial \psi}{\partial x_2} + \delta \frac{\partial \psi}{\partial x_3} + \dots \right\}. \end{aligned} \quad (20.66)$$

According to Eq. (20.65) the operator of the momentum of the state $|\psi\rangle$ (let's call it the *total* momentum) is equal to

$$\hat{p}_{\text{total}} = \frac{\hbar}{i} \left\{ \frac{\partial}{\partial x_1} + \frac{\partial}{\partial x_2} + \frac{\partial}{\partial x_3} + \dots \right\}.$$

But this is just the same as

$$\hat{p}_{\text{total}} = \hat{p}_{x1} + \hat{p}_{x2} + \hat{p}_{x3} + \dots. \quad (20.67)$$

The operators of momentum obey the rule that the total momentum is the sum of the momenta of all the parts. Everything holds together nicely, and many of the things we have been saying are consistent with each other.

20-6 Angular momentum

Let's for fun look at another operation—the operation of orbital angular momentum. In Chapter 17 we defined an operator \hat{J}_z in terms of $\hat{R}_z(\varphi)$, the operator of a rotation by the angle φ about the z -axis. We consider here a system described simply by a single wave function $\psi(\mathbf{r})$, which is a function of coordinates only, and does not take into account the fact that the electron may have its spin either up or down. That is, we want for the moment to disregard *intrinsic* angular momentum and think about only the *orbital* part. To keep the distinction clear, we'll call the orbital operator \hat{L}_z , and define it in terms of the operator of a rotation by an infinitesimal angle ϵ by

$$\hat{R}_z(\epsilon) |\psi\rangle = \left(1 + \frac{i}{\hbar} \epsilon \hat{L}_z\right) |\psi\rangle.$$

(Remember, this definition applies only to a state $|\psi\rangle$ which has no internal spin variables, but depends only on the coordinates $\mathbf{r} = x, y, z$.) If we look at the state $|\psi\rangle$ in a new coordinate system, rotated about the z -axis by the small angle ϵ , we see a new state

$$|\psi'\rangle = \hat{R}_z(\epsilon) |\psi\rangle.$$

If we choose to describe the state $|\psi\rangle$ in the coordinate representation—that is, by its wave function $\psi(\mathbf{r})$, we would expect to be able to write

$$\psi'(\mathbf{r}) = \left(1 + \frac{i}{\hbar} \epsilon \hat{L}_z\right) \psi(\mathbf{r}). \quad (20.68)$$

What is \hat{L}_z ? Well, a point P at x and y in the *new* coordinate system (really x' and y' , but we will drop the primes) was formerly at $x - \epsilon y$ and $y + \epsilon x$, as you can see from Fig. 20-2. Since the amplitude for the electron to be at P isn't changed by the rotation of the coordinates we can write

$$\psi'(x, y, z) = \psi(x + \epsilon y, y - \epsilon x, z) = \psi(x, y, z) + \epsilon y \frac{\partial \psi}{\partial x} - \epsilon x \frac{\partial \psi}{\partial y}$$

(remembering that ϵ is a small angle). This means that

$$\hat{L}_z = \frac{\hbar}{i} \left(x \frac{\partial}{\partial y} - y \frac{\partial}{\partial x} \right). \quad (20.69)$$

That's our answer. But notice. It is equivalent to

$$\hat{L}_z = x \hat{p}_y - y \hat{p}_x. \quad (20.70)$$

Returning to our quantum-mechanical operators, we can write

$$\hat{L}_z = x \hat{p}_y - y \hat{p}_x. \quad (20.71)$$

This formula is easy to remember because it looks like the familiar formula of classical mechanics; it is the z -component of

$$\mathbf{L} = \mathbf{r} \times \mathbf{p}. \quad (20.72)$$

One of the fun parts of this operator business is that many classical equations get carried over into a quantum-mechanical form. Which ones don't? There had better be some that don't come out right, because if everything did, then there would be nothing different about quantum mechanics. There would be no new physics. Here is one equation which is different. In classical physics

$$xp_x - p_x x = 0.$$

What is it in quantum mechanics?

$$\hat{x}\hat{p}_x - \hat{p}_x\hat{x} = ?$$

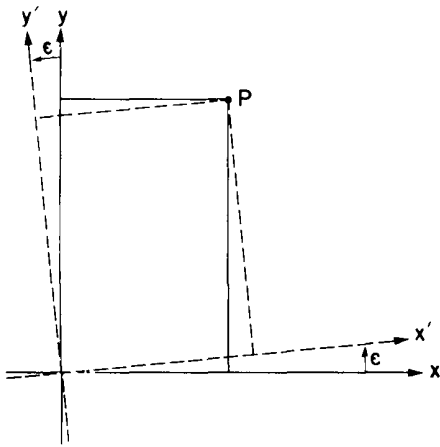


Fig. 20-2. Rotation of the axes around the z -axis by the small angle ϵ .

Let's work it out in the x -representation. So that we'll know what we are doing we put in some wave function $\psi(x)$. We have

$$x\mathcal{P}_x\psi(x) - \hat{\mathcal{P}}_xx\psi(x),$$

or

$$x \frac{\hbar}{i} \frac{\partial}{\partial x} \psi(x) - \frac{\hbar}{i} \frac{\partial}{\partial x} x\psi(x).$$

Remember now that the derivatives operate on everything to the right. We get

$$x \frac{\hbar}{i} \frac{\partial \psi}{\partial x} - \frac{\hbar}{i} \psi(x) - \frac{\hbar}{i} x \frac{\partial \psi}{\partial x} = -\frac{\hbar}{i} \psi(x). \quad (20.73)$$

The answer is *not* zero. The whole operation is equivalent simply to multiplication by $-\hbar/i$:

$$\hat{x}\hat{p}_x - \hat{p}_x\hat{x} = -\frac{\hbar}{i}. \quad (20.74)$$

If Plank's constant were zero, the classical and quantum results would be the same, and there would be no quantum mechanics to learn!

Incidentally, if any two operators \hat{A} and \hat{B} , when taken together like this:

$$\hat{A}\hat{B} - \hat{B}\hat{A},$$

do *not* give zero, we say that "the operators do not commute." And an equation such as (20.74) is called a "commutation rule." You can see that the commutation rule for p_x and y is

$$\hat{p}_x\hat{y} - \hat{y}\hat{p}_x = 0.$$

There is another very important commutation rule that has to do with angular momenta. It is

$$\hat{L}_x\hat{L}_y - \hat{L}_y\hat{L}_x = i\hbar\hat{L}_z. \quad (20.75)$$

You can get some practice with \hat{x} and \hat{p} operators by proving it for yourself.

It is interesting to notice that operators which do not commute can also occur in classical physics. We have already seen this when we have talked about rotation in space. If you rotate something, such as a book, by 90° around x and then 90° around y , you get something different from rotating first by 90° around y and then by 90° around x . It is, in fact, just this property of space that is responsible for Eq. (20.75).

20-7 The change of averages with time

Now we want to show you something else. How do averages change with time? Suppose for the moment that we have an operator \hat{A} , which does not itself have time in it in any obvious way. We mean an operator like \hat{x} or \hat{p} . (We exclude things like, say, the operator of some external potential that was being varied with time, such as $V(x, t)$.) Now suppose we calculate $\langle A \rangle_{av}$, in some state $|\psi\rangle$, which is

$$\langle A \rangle_{av} = \langle \psi | \hat{A} | \psi \rangle. \quad (20.76)$$

How will $\langle A \rangle_{av}$ depend on time? Why should it? One reason might be that the operator itself depended explicitly on time—for instance, if it had to do with a time-varying potential like $V(x, t)$. But even if the operator does not depend on t , say, for example, the operator $\hat{A} = \hat{x}$, the corresponding average may depend on time. Certainly the average position of a particle could be moving. How does such a motion come out of Eq. (20.76) if \hat{A} has no time dependence? Well, the state $|\psi\rangle$ might be changing with time. For nonstationary states we have often shown a time dependence explicitly by writing a state as $|\psi(t)\rangle$. We want to show that the rate of change of $\langle A \rangle_{av}$ is given by a new operator we will call $\dot{\hat{A}}$. Remember that \hat{A} is an operator, so that putting a dot over the A does not here mean taking

the time derivative, but is just a way of writing a *new* operator \hat{A} which is defined by

$$\frac{d}{dt} \langle A \rangle_{av} = \langle \psi | \hat{A} | \psi \rangle. \quad (20.77)$$

Our problem is to find the operator \hat{A} .

First, we know that the rate of change of a state is given by the Hamiltonian. Specifically,

$$i\hbar \frac{d}{dt} | \psi(t) \rangle = \hat{H} | \psi(t) \rangle. \quad (20.78)$$

This is just the abstract way of writing our original definition of the Hamiltonian:

$$i\hbar \frac{dC_i}{dt} = \sum_{ij} H_{ij} C_j. \quad (20.79)$$

If we take the complex conjugate of this equation, it is equivalent to

$$-i\hbar \frac{d}{dt} \langle \psi(t) | = \langle \psi(t) | \hat{H}. \quad (20.80)$$

Next, see what happens if we take the derivatives with respect to t of Eq. (20.76). Since each ψ depends on t , we have

$$\frac{d}{dt} \langle A \rangle_{av} = \left(\frac{d}{dt} \langle \psi | \right) \hat{A} | \psi + \langle \psi | \hat{A} \left(\frac{d}{dt} | \psi \rangle \right). \quad (20.81)$$

Finally, using the two equations in (20.78) and (20.79) to replace the derivatives, we get

$$\frac{d}{dt} \langle A \rangle_{av} = \frac{i}{\hbar} \{ \langle \psi | \hat{H} \hat{A} | \psi \rangle - \langle \psi | \hat{A} \hat{H} | \psi \rangle \}.$$

This equation is the same as

$$\frac{d}{dt} \langle A \rangle_{av} = \frac{i}{\hbar} \langle \psi | (\hat{H} \hat{A} - \hat{A} \hat{H}) | \psi \rangle.$$

Comparing this equation with Eq. (20.77), you see that

$$\hat{A} = \frac{i}{\hbar} (\hat{H} \hat{A} - \hat{A} \hat{H}). \quad (20.82)$$

That is our interesting proposition, and it is true for any operator \hat{A} .

Incidentally, if the operator \hat{A} should *itself* be time dependent, we would have had

$$\hat{A} = \frac{i}{\hbar} (\hat{H} \hat{A} - \hat{A} \hat{H}) + \frac{\partial \hat{A}}{\partial t}. \quad (20.83)$$

Let us try out Eq. (20.82) on some example to see whether it really makes sense. For instance, what operator corresponds to \hat{x} ? We say it should be

$$\hat{x} = \frac{i}{\hbar} (\hat{H} \hat{x} - \hat{x} \hat{H}). \quad (20.84)$$

What is this? One way to find out is to work it through in the coordinate representation using the algebraic operator for $\hat{\mathcal{H}}$. In this representation the commutator is

$$\hat{\mathcal{H}}x - x\hat{\mathcal{H}} = \left\{ \frac{\hbar^2}{2m} \frac{d^2}{dx^2} + V(x) \right\} x - x \left\{ \frac{\hbar^2}{2m} \frac{d^2}{dx^2} + V(x) \right\}.$$

If you operate with this on any wave function $\psi(x)$ and work out all of the derivatives where you can, you end up after a little work with

$$-\frac{\hbar^2}{2m} \frac{d\psi}{dx}.$$

But this is just the same as

$$-i \frac{\hbar}{m} \hat{p}_x \psi,$$

so we find that

$$\hat{H}\hat{x} = \hat{x}\hat{H} = -i \frac{\hbar}{m} p_x \quad (20.85)$$

or that

$$\dot{\bar{x}} = \frac{\bar{p}_x}{m}. \quad (20.86)$$

A pretty result. It means that if the mean value of x is changing with time the drift of the center of gravity is the same as the mean momentum divided by m . Exactly like classical mechanics.

Another example. What is the rate of change of the average momentum of a state? Same game. Its operator is

$$\hat{p} = \frac{i}{\hbar} (\hat{H}\hat{p} - \hat{p}\hat{H}). \quad (20.87)$$

Again you can work it out in the x representation. Remember that \hat{p} becomes d/dx , and this means that you will be taking the derivative of the potential energy V (in the $\hat{\mathcal{C}}$)—but only in the second term. It turns out that it is the only term which does not cancel, and you find that

$$\hat{\mathcal{C}}\hat{p} - \hat{p}\hat{\mathcal{C}} = -i\hbar \frac{dV}{dx}$$

or that

$$\hat{p} = - \frac{dV}{dx}. \quad (20.88)$$

Again the classical result. The right-hand side is the force, so we have derived Newton's law! But remember—these are the laws for the *operators* which give the *average* quantities. They do not describe what goes on in detail inside an atom.

Quantum mechanics has the essential difference that $\hat{p}\hat{x}$ is not equal to $\hat{x}\hat{p}$. They differ by a little bit—by the small number \hbar . But the whole wondrous complications of interference, waves, and all, result from the little fact that $\hat{x}\hat{p} - \hat{p}\hat{x}$ is not quite zero.

The history of this idea is also interesting. Within a period of a few months in 1926, Heisenberg and Schrödinger independently found correct laws to describe atomic mechanics. Schrödinger invented his wave function $\psi(x)$ and found his equation. Heisenberg, on the other hand, found that nature could be described by classical equations, except that $xp - px$ should be equal to \hbar/i , which he could make happen by defining them in terms of special kinds of matrices. In our language he was using the energy-representation, with its matrices. Both Heisenberg's matrix algebra and Schrödinger's differential equation explained the hydrogen atom. A few months later Schrödinger was able to show that the two theories were equivalent—as we have seen here. But the two different mathematical forms of quantum mechanics were discovered independently.

The Schrödinger Equation in a Classical Context: A Seminar on Superconductivity

21-1 Schrödinger's equation in a magnetic field

This lecture is only for entertainment. I would like to give the lecture in a somewhat different style—just to see how it works out. It's not a part of the course—in the sense that it is not supposed to be a last minute effort to teach you something new. But, rather, I imagine that I'm giving a seminar or research report on the subject to a more advanced audience, to people who have already been educated in quantum mechanics. The main difference between a seminar and a regular lecture is that the seminar speaker does not carry out all the steps, or all the algebra. He says: "If you do such and such, this is what comes out," instead of showing all of the details. So in this lecture I'll describe the ideas all the way along but just give you the *results* of the computations. You should realize that you're not supposed to understand everything immediately, but believe (more or less) that things would come out if you went through the steps.

All that aside, this is a subject I *want* to talk about. It is recent and modern and would be a perfectly legitimate talk to give at a research seminar. My subject is the Schrödinger equation in a classical setting—the case of superconductivity.

Ordinarily, the wave function which appears in the Schrödinger equation applies to only one or two particles. And the wave function itself is not something that has a classical meaning—unlike the electric field, or the vector potential, or things of that kind. The wave function for a single particle *is* a "field"—in the sense that it is a function of position—but it does not generally have a classical significance. Nevertheless, there are some situations in which a quantum mechanical wave function *does* have classical significance, and they are the ones I would like to take up. The peculiar quantum mechanical behavior of matter on a small scale doesn't usually make itself felt on a large scale except in the standard way that it produces Newton's laws—the laws of the so-called classical mechanics. But there are certain situations in which the peculiarities of quantum mechanics can come out in a special way on a large scale.

At low temperatures, when the energy of a system has been reduced very, very low, instead of a large number of states being involved, only a very, very small number of states near the ground state are involved. Under those circumstances the quantum mechanical character of that ground state can appear on a macroscopic scale. It is the purpose of this lecture to show a connection between quantum mechanics and large-scale effects—not the usual discussion of the way that quantum mechanics reproduces Newtonian mechanics on the average, but a special situation in which quantum mechanics will produce its own characteristic effects on a large or "macroscopic" scale.

I will begin by reminding you of some of the properties of the Schrödinger equation.† I want to describe the behavior of a particle in a magnetic field using the Schrödinger equation, because the superconductive phenomena are involved with magnetic fields. An external magnetic field is described by a vector potential, and the problem is: what are the laws of quantum mechanics in a vector potential? The principle that describes the behavior of quantum mechanics in a vector potential is very simple. The amplitude that a particle goes from one place to another along a certain route when there's a field present is the same as the ampli-

† I'm not really reminding you, because I haven't shown you some of these equations before; but remember the spirit of this seminar.

- 21-1 Schrödinger's equation in a magnetic field**
- 21-2 The equation of continuity for probabilities**
- 21-3 Two kinds of momentum**
- 21-4 The meaning of the wave function**
- 21-5 Superconductivity**
- 21-6 The Meissner effect**
- 21-7 Flux quantization**
- 21-8 The dynamics of superconductivity**
- 21-9 The Josephson junction**

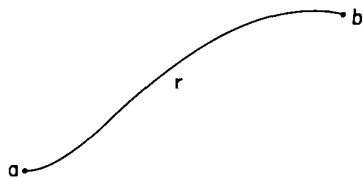


Fig. 21-1. The amplitude to go from a to b along the path Γ is proportional to $\exp(iq/\hbar) \int_a^b \mathbf{A} \cdot d\mathbf{s}$.

tude that it would go along the same route when there's no field, multiplied by the exponential of the line integral of the vector potential, times the electric charge divided by Planck's constant¹ (see Fig. 21-1):

$$\langle b | a \rangle_{\mathbf{A}} = \langle b | a \rangle_{\mathbf{A}=0} \cdot \exp \left\{ \frac{iq}{\hbar} \int_a^b \mathbf{A} \cdot d\mathbf{s} \right\}. \quad (21.1)$$

It is a basic statement of quantum mechanics.

Now without the vector potential the Schrodinger equation of a charged particle (nonrelativistic, no spin) is

$$-\frac{\hbar}{i} \frac{\partial \psi}{\partial t} = \hat{\mathcal{C}}\psi = \frac{1}{2m} \left(\frac{\hbar}{i} \nabla \right) \cdot \left(\frac{\hbar}{i} \nabla \right) \psi + q\phi\psi, \quad (21.2)$$

where ϕ is the electric potential so that $q\phi$ is the potential energy.[†] Equation (21.1) is equivalent to the statement that in a magnetic field the gradients in the Hamiltonian are replaced in each case by the gradient minus $q\mathbf{A}$, so that Eq. (21.2) becomes

$$-\frac{\hbar}{i} \frac{\partial \psi}{\partial t} = \hat{\mathcal{C}}\psi = \frac{1}{2m} \left(\frac{\hbar}{i} \nabla - q\mathbf{A} \right) \cdot \left(\frac{\hbar}{i} \nabla - q\mathbf{A} \right) \psi + q\phi\psi. \quad (21.3)$$

This is the Schrodinger equation for a particle with charge q moving in an electromagnetic field \mathbf{A} , ϕ (nonrelativistic, no spin).

To show that this is true I'd like to illustrate by a simple example in which instead of having a continuous situation we have a line of atoms along the x -axis with the spacing b and we have an amplitude $-K$ for an electron to jump from one atom to another when there is no field.[‡] Now according to Eq. (21.1) if there's a vector potential in the x -direction $A_x(x, t)$, the amplitude to jump will be altered from what it was before by a factor $\exp(iq/\hbar A_x b)$, the exponent being iq/\hbar times the vector potential integrated from one atom to the next. For simplicity we will write $(q/\hbar)A_x \equiv f(x)$, since A_x will, in general, depend on x . If the amplitude to find the electron at the atom " n " located at x is called $C(x) \equiv C_n$, then the rate of change of that amplitude is given by the following equation:

$$\begin{aligned} -\frac{\hbar}{i} \frac{\partial}{\partial t} C(x) = & E_0 C(x) - K e^{-ibf(x+b/2)} C(x+b) \\ & - K e^{+ibf(x-b/2)} C(x-b). \end{aligned} \quad (21.4)$$

There are three pieces. First, there's some energy E_0 if the electron is located at x . As usual, that gives the term $E_0 C(x)$. Next, there is the term $-KC(x+b)$, which is the amplitude for the electron to have jumped backwards one step from atom " $n+1$," located at $x+b$. However, in doing so in a vector potential, the phase of the amplitude must be shifted according to the rule in Eq. (21.1). If A_x is not changing appreciably in one atomic spacing, the integral can be written as just the value of A_x at the midpoint, times the spacing b . So (iq/\hbar) times the integral is just $b f(x+b/2)$. Since the electron is jumping backwards, I showed this phase shift with a minus sign. That gives the second piece. In the same manner there's a certain amplitude to have jumped from the other side, but this time we need the vector potential at a distance $(b/2)$ on the other side of x , times the distance b . That gives the third piece. The sum gives the equation for the amplitude to be at x in a vector potential.

Now we know that if the function $C(x)$ is smooth enough (we take the long wavelength limit), and if we let the atoms get closer together, Eq. (16.4) will approach the behavior of an electron in free space. So the next step is to expand both sides of (21.4) in powers of b , assuming b is very small. For example, if b is zero the right-hand side is just $(E_0 - 2K)C(x)$, so in the zeroth approximation

¹ Volume II, Section 15-5.

[†] Not to be confused with our earlier use of ϕ for a state label!

[‡] K is the same quantity that was called A in the problem of a linear lattice with no magnetic field. See Chapter 13.

the energy is $E_0 - 2K$. Next comes the terms in b . But because the two exponentials have opposite signs, only even powers of b remain. So if you make a Taylor expansion of $C(x)$, of $f(x)$, and of the exponentials, and then collect the terms in b^2 , you get

$$-\frac{\hbar}{i} \frac{\partial C(x)}{\partial t} = E_0 C(x) - 2KC(x) - Kb^2 \{C''(x) - 2if(x)C'(x) - if'(x)C(x) - f^2(x)C(x)\}. \quad (21.5)$$

(The “primes” mean differentiation with respect to x .)

Now this horrible combination of things looks quite complicated. But mathematically it's exactly the same as

$$-\frac{\hbar}{i} \frac{\partial C(x)}{\partial t} = (E_0 - 2K)C(x) - Kb^2 \left[\frac{\partial}{\partial x} - if(x) \right] \left[\frac{\partial}{\partial x} - if(x) \right] C(x). \quad (21.6)$$

The second bracket operating on $C(x)$ gives $C'(x)$ plus $if(x)C(x)$. The first bracket operating on these two terms gives the C'' term and terms in the first derivative of $f(x)$ and the first derivative of $C(x)$. Now remember that the solutions for zero magnetic field² represent a particle with an effective mass m_{eff} given by

$$Kb^2 = \frac{\hbar}{m_{\text{eff}}}.$$

If you then set $E_0 = -2K$, and put back $f(x) = (q/\hbar)A_x$, you can easily check that Eq. (21.6) is the same as the first part of Eq. (21.3). (The origin of the potential energy term is well known, so I haven't bothered to include it in this discussion.) The proposition of Eq. (21.1) that the vector potential changes all the amplitudes by the exponential factor is the same as the rule that the momentum operator, $(\hbar/i)\nabla$ gets replaced by

$$\frac{\hbar}{i} \nabla - qA,$$

as you see in the Schrödinger equation of (21.3).

21-2 The equation of continuity for probabilities

Now I turn to a second point. An important part of the Schrödinger equation for a single particle is the idea that the probability to find the particle at a position is given by the absolute square of the wave function. It is also characteristic of the quantum mechanics that probability is conserved in a local sense. When the probability of finding the electron somewhere decreases, while the probability of the electron being elsewhere increases (keeping the total probability unchanged), something must be going on in between. In other words, the electron has a continuity in the sense that if the probability decreases at one place and builds up at another place, there must be some kind of flow between. If you put a wall, for example, in the way, it will have an influence and the probabilities will not be the same. So the conservation of probability alone is not the complete statement of the conservation law, just as the conservation of energy alone is not as deep and important as the *local* conservation of energy.³ If energy is disappearing, there must be a flow of energy to correspond. In the same way, we would like to find a “current” of probability such that if there is any change in the probability density (the probability of being found in a unit volume), it can be considered as coming from an inflow or an outflow due to some current. This current would be a vector which could be interpreted this way—the x component would be the net probability per second and per unit area that a particle passes in the x direction across a plane parallel to the y - z plane. Passage toward $+x$ is considered a positive flow, and passage in the opposite direction, a negative flow.

² Section 13-3.

³ Volume II, Section 27-1.

Is there such a current? Well, you know that the probability density $P(\mathbf{r}, t)$ is given in terms of the wave function by

$$P(\mathbf{r}, t) = \psi^*(\mathbf{r}, t)\psi(\mathbf{r}, t). \quad (21.7)$$

I am asking: Is there a current \mathbf{J} such that

$$\frac{\partial P}{\partial t} = -\nabla \cdot \mathbf{J}? \quad (21.8)$$

If I take the time derivative of Eq. (21.7), I get two terms:

$$\frac{\partial P}{\partial t} = \psi^* \frac{\partial \psi}{\partial t} + \psi \frac{\partial \psi^*}{\partial t}. \quad (21.9)$$

Now use the Schrödinger equation—Eq. (21.3)—for $\partial\psi/\partial t$; and take the complex conjugate of it to get $\partial\psi^*/\partial t$ —each i gets its sign reversed. You get

$$\begin{aligned} \frac{\partial P}{\partial t} = & -\frac{i}{\hbar} \psi^* \frac{1}{2m} \left(\frac{\hbar}{i} \nabla - q\mathbf{A} \right) \cdot \left(\frac{\hbar}{i} \nabla - q\mathbf{A} \right) \psi + e\phi \psi^* \psi \\ & - \psi \frac{1}{2m} \left(\frac{\hbar}{i} \nabla + q\mathbf{A} \right) \cdot \left(\frac{\hbar}{i} \nabla + q\mathbf{A} \right) \psi^* - e\phi \psi \psi^*. \end{aligned} \quad (21.10)$$

The potential terms and a lot of other stuff cancel out. And it turns out that what is left can indeed be written as a perfect divergence. The whole equation is equivalent to

$$\frac{\partial P}{\partial t} = -\nabla \cdot \left\{ \frac{1}{2m} \psi^* \left(\frac{\hbar}{i} \nabla - q\mathbf{A} \right) \psi + \psi \left(-\frac{\hbar}{i} \nabla - q\mathbf{A} \right) \psi^* \right\}. \quad (21.11)$$

It is really not as complicated as it seems. It is a symmetrical combination of ψ^* times a certain operation on ψ , plus ψ^* times the complex conjugate operation on ψ . It is some quantity plus its own complex conjugate, so the whole thing is real—as it ought to be. The operation can be remembered this way: it is just the momentum operator \hat{p} minus $q\mathbf{A}$. I could write the current in Eq. (21.8) as

$$\mathbf{J} = \frac{1}{2} \left\{ \left[\frac{\hat{p} - q\mathbf{A}}{m} \psi \right]^* \psi + \psi^* \left[\frac{\hat{p} - q\mathbf{A}}{m} \right] \psi \right\}. \quad (21.12)$$

There is then a current \mathbf{J} which completes Eq. (21.8).

Equation (21.10) shows that the probability is conserved locally. If a particle disappears from one region it cannot appear in another without something going on in between. Imagine that the first region is surrounded by a closed surface far enough out that there is zero probability to find the electron at the surface. The total probability to find the electron somewhere inside the surface is the volume integral of P . But according to Gauss's theorem the volume integral of the divergence \mathbf{J} is equal to the surface integral of \mathbf{J} . If ψ is zero at the surface, Eq. (21.10) says that \mathbf{J} is zero, so the total probability to find the particle inside can't change. Only if some of the probability approaches the boundary can some of it leak out. We can say that it only gets out by moving through the surface—and that is local conservation.

21-3 Two kinds of momentum

The equation for the current is rather interesting, and sometimes causes a certain amount of worry. You would think the current would be something like the density of particles times the velocity. The density should be something like $\psi\psi^*$, which is o.k. And each term in Eq. (21.12) looks like the typical form for the average-value of the operator

$$\frac{\hat{p} - q\mathbf{A}}{m}, \quad (21.13)$$

so maybe we should think of it as the velocity of flow. It looks as though we have two suggestions for relations of velocity to momentum, because we would also think that momentum divided by mass, $\hat{\Phi}/m$, should be a velocity. The two possibilities differ by the vector potential.

It happens that these two possibilities were also discovered in classical physics, when it was found that momentum could be defined in two ways.⁴ One of them is called “kinematic momentum,” but for absolute clarity I will in this lecture call it the “ mv -momentum.” This is the momentum obtained by multiplying mass by velocity. The other is a more mathematical, more abstract momentum, sometimes called the “dynamical momentum,” which I’ll call “ p -momentum.” The two possibilities are

$$mv\text{-momentum} = mv, \quad (21.14)$$

$$p\text{-momentum} = mv + qA. \quad (21.15)$$

It turns out that in quantum mechanics with magnetic fields it is the p -momentum which is connected to the gradient operator $\hat{\Phi}$, so it follows that (21.13) is the operator of a velocity.

I’d like to make a brief digression to show you what this is all about—why there must be something like Eq. (21.15) in the quantum mechanics. The wave function changes with time according to the Schrödinger equation in Eq. (21.3). If I would suddenly change the vector potential, the wave function wouldn’t change at the first instant; only its rate of change changes. Now think of what would happen in the following circumstance. Suppose I have a long solenoid, in which I can produce a flux of magnetic field (B -field), as shown in Fig. 21-2. And there is a charged particle sitting nearby. Suppose this flux nearly instantaneously builds up from zero to something. I start with zero vector potential and then I turn on a vector potential. That means that I produce suddenly a circumferential vector potential A . You’ll remember that the line integral of A around a loop is the same as the flux of B through the loop.⁵ Now what happens if I suddenly turn on a vector potential? According to the quantum mechanical equation the sudden change of A does not make a sudden change of ψ ; the wave function is still the same. So the gradient is also unchanged.

But remember what happens electrically when I suddenly turn on a flux. During the short time that the flux is rising, there’s an electric field generated whose line integral is the rate of change of the flux with time:

$$E = - \frac{\partial A}{\partial t}. \quad (21.16)$$

That electric field is enormous if the flux is changing rapidly, and it gives a force on the particle. The force is the charge times the electric field, and so during the build up of the flux the particle obtains a total impulse (that is, a change in mv) equal to $-qA$. In other words, if you suddenly turn on a vector potential at a charge, this charge immediately picks up an “ mv ” momentum equal to $-qA$. But there is something that isn’t changed immediately and that’s the difference between mv and $-qA$. And so the sum $p = mv + qA$ is something which is not changed when you make a sudden change in the vector potential. This quantity p is what we have called the p -momentum and is of importance in classical mechanics in the theory of dynamics, but it also has a direct significance in quantum mechanics. It depends on the character of the wave function, and it is the one to be identified with the operator

$$\hat{\Phi} = \frac{\hbar}{i} \nabla.$$

⁴ See, for example, J. D. Jackson, *Classical Electrodynamics*, John Wiley and Sons, Inc. New York (1962), p. 408.

⁵ Volume II, Chapter 14, Section 14-1.

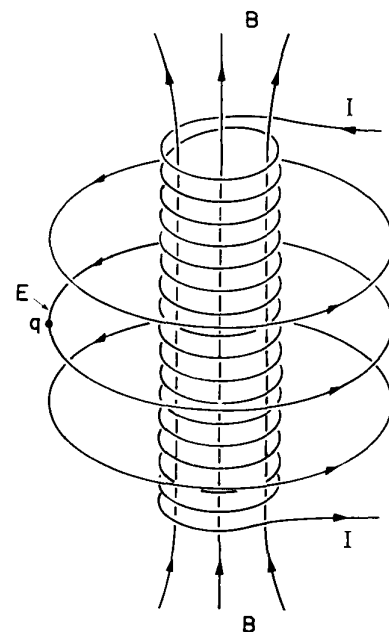


Fig. 21-2. The electric field outside a solenoid with an increasing current.

21-4 The meaning of the wave function

When Schrodinger first discovered his equation he discovered the conservation law of Eq. (21.9) as a consequence of his equation. But he imagined incorrectly that P was the electric charge density of the electron and that \mathbf{J} was the electric current density, so he thought that the electrons interacted with the electromagnetic field through these charges and currents. When he solved his equations for the hydrogen atom and calculated ψ , he wasn't calculating the probability of anything—there were no amplitudes at that time—the interpretation was completely different. The atomic nucleus was stationary but there were currents moving around; the charges P and currents \mathbf{J} would generate electromagnetic fields and the thing would radiate light. He soon found on doing a number of problems that it didn't work out quite right. It was at this point that Born made an essential contribution to our ideas regarding quantum mechanics. It was Born who correctly (as far as we know) interpreted the ψ of the Schrodinger equation in terms of a probability amplitude—that very difficult idea that the square of the amplitude is not the charge density but is only the probability per unit volume of finding an electron there, and that when you do find the electron some place the entire charge is there. That whole idea is due to Born.

The wave function $\psi(\mathbf{r})$ for an electron in an atom does not, then, describe a smeared-out electron with a smooth charge density. The electron is either here, or there, or somewhere else, but wherever it is, it is a point charge. On the other hand, think of a situation in which there are an enormous number of particles in exactly the same state, a very large number of them with exactly the same wave function. Then what? One of them is here and one of them is there, and the probability of finding any one of them at a given place is proportional to $\psi\psi^*$. But since there are so many particles, if I look in any volume $dx\,dy\,dz$ I will generally find a number close to $\psi\psi^* dx\,dy\,dz$. So in a situation in which ψ is the wave function for each of an enormous number of particles which are all in the same state, $\psi\psi^*$ can be interpreted as the density of particles. If, under these circumstances, each particle carries the same charge q , we can, in fact, go further and interpret $\psi^*\psi$ as the density of *electricity*. Normally, $\psi\psi^*$ is given the dimensions of a probability density, then ψ should be multiplied by q to give the dimensions of a charge density. For our present purposes we can put this constant factor into ψ , and take $\psi\psi^*$ itself as the electric charge density. With this understanding, \mathbf{J} (the current of probability I have calculated) becomes directly the electric current density.

So in the situation in which we can have very many particles in exactly the same state, there is possible a new physical interpretation of the wave functions. The charge density and the electric current can be calculated directly from the wave functions and the wave functions take on a physical meaning which extends into classical, macroscopic situations.

Something similar can happen with neutral particles. When we have the wave function of a single photon, it is the amplitude to find a photon somewhere. Although we haven't ever written it down there is an equation for the photon wave function analogous to the Schrödinger equation for the electron. The photon equation is just the same as Maxwell's equations for the electromagnetic field, and the wave function is the same as the vector potential \mathbf{A} . The wave function turns out to be just the vector potential. The quantum physics is the same thing as the classical physics because photons are noninteracting Bose particles and many of them can be in the same state—as you know, they *like* to be in the same state. The moment that you have billions in the same state (that is, in the same electromagnetic wave), you can measure the wave function, which is the vector potential, directly. Of course, it worked historically the other way. The first observations were on situations with many photons in the same state, and so we were able to discover the correct equation for a single photon by observing directly with our hands on a macroscopic level the nature of wave function.

Now the trouble with the electron is that you cannot put more than one in the same state. Therefore, it was long believed that the wave function of the

Schrodinger equation would never have a macroscopic representation analogous to the macroscopic representation of the amplitude for photons. On the other hand, it is now realized that the phenomena of superconductivity presents us with just this situation.

21-5 Superconductivity

As you know, very many metals become superconducting below a certain temperature⁶—the temperature is different for different metals. When you reduce the temperature sufficiently the metals conduct electricity without any resistance. This phenomenon has been observed for a very large number of metals but not for all, and the theory of this phenomenon has caused a great deal of difficulty. It took a very long time to understand what was going on inside of superconductors, and I will only describe enough of it for our present purposes. It turns out that due to the interactions of the electrons with the vibrations of the atoms in the lattice, there is a small net effective *attraction* between the electrons. The result is that the electrons form together, if I may speak very qualitatively and crudely, bound pairs.

Now you know that a single electron is a Fermi particle. But a bound pair would act as a Bose particle, because if I exchange both electrons in a pair I change the sign of the wave function twice, and that means that I don't change anything. A pair *is* a Bose particle.

The energy of pairing—that is, the net attraction—is very, very weak. Only a tiny temperature is needed to throw the electrons apart by thermal agitation, and convert them back to “normal” electrons. But when you make the temperature sufficiently low that they have to do their very best to get into the absolutely lowest state; then they do collect in pairs.

I don't wish you to imagine that the pairs are really held together very closely like a point particle. As a matter of fact, one of the great difficulties of understanding this phenomena originally was that that is not the way things are. The two electrons which form the pair are really spread over a considerable distance; and the mean distance between pairs is relatively smaller than the size of a single pair. Several pairs are occupying the same space at the same time. Both the reason why electrons in a metal form pairs and an estimate of the energy given up in forming a pair have been a triumph of recent times. This fundamental point in the theory of superconductivity was first explained in the theory of Bardeen, Cooper, and Schrieffer,⁷ but that is not the subject of this seminar. We will accept, however, the idea that the electrons do, in some manner or other, work in pairs, that we can think of these pairs as behaving more or less like particles, and that we can therefore talk about the wave function for a “pair.”

Now the Schrödinger equation for the pair will be more or less like Eq. (21.3). There will be one difference in that the charge q will be twice the charge of an electron. Also, we don't know the inertia—or effective mass—for the pair in the crystal lattice, so we don't know what number to put in for m . Nor should we think that if we go to very high frequencies (or short wavelengths), this is exactly the right form, because the kinetic energy that corresponds to very rapidly varying wave functions may be so great as to break up the pairs. At finite temperatures there are always a few pairs which are broken up according to the usual Boltzmann theory. The probability that a pair is broken is proportional to $\exp(-E_{\text{pair}}/kT)$. The electrons that are not bound in pairs are called “normal” electrons and will move around in the crystal in the ordinary way. I will, however, consider only the situation at essentially zero temperature—or, in any case, I will disregard the complications produced by those electrons which are not in pairs.

⁶ First discovered by Onnes in 1911; H. K. Onnes, Comm. Phys. Lab., Univ. Leyden, Nos. 119, 120, 122 (1911). You will find a nice up-to-date discussion of the subject in E. A. Lynton, *Superconductivity*, John Wiley and Sons, Inc., New York, 1962.

⁷ J. Bardeen, L. N. Cooper, and J. R. Schrieffer, *Phys. Rev.* **108**, 1175 (1957).

Since electron pairs are bosons, when there are a lot of them in a given state there is an especially large amplitude for other pairs to go to the same state. So nearly all of the pairs will be locked down at the lowest energy in *exactly the same state*—it won't be easy to get one of them into another state. There's more amplitude to go into the same state than into an unoccupied state by the famous factor \sqrt{n} , where n is the occupancy of the lowest state. So we would expect all the pairs to be moving in the same state.

What then will our theory look like? I'll call ψ the wave function of a pair in the lowest energy state. However, since $\psi\psi^*$ is going to be proportional to the charge density ρ , I can just as well write ψ as the square root of the charge density times some phase factor:

$$\psi(r) = \rho(r)e^{i\theta(r)}, \quad (21.17)$$

where ρ and θ are real functions of r . (Any complex function can, of course, be written this way.) It's clear what we mean when we talk about the charge density, but what is the physical meaning of the phase θ of the wave function? Well, let's see what happens if we substitute $\psi(r)$ into Eq. (21.12), and express the current density in terms of these new variables ρ and θ . It's just a change of variables and I won't go through all the algebra, but it comes out

$$\mathbf{J} = \frac{\hbar}{m} \left(\nabla\theta - \frac{q}{\hbar} \mathbf{A} \right) \rho. \quad (21.18)$$

Since both the current density and the charge density have a direct physical meaning for the superconducting electron gas, both ρ and θ are real things. The phase is just as observable as ρ ; it is a piece of the current density \mathbf{J} . The *absolute* phase is not observable, but if the gradient of the phase is known everywhere, the phase is known except for a constant. You can define the phase at one point, and then the phase everywhere is determined.

Incidentally, the equation for the current can be analyzed a little nicer, when you think that the current density \mathbf{J} is *in fact* the charge density times the velocity of motion of the fluid of electrons, or $\rho\mathbf{v}$. Equation (21.18) is then equivalent to

$$m\mathbf{v} = \hbar \nabla\theta - q\mathbf{A}. \quad (21.19)$$

Notice that there are two pieces in the $m\mathbf{v}$ -momentum; one is a contribution from the vector potential, and the other, a contribution from the behavior of the wave function. In other words, the quantity $\hbar \nabla\theta$ is just what we have called the p -momentum.

21-6 The Meissner effect

Now we can describe some of the phenomena of superconductivity. First, there is no electrical resistance. There's no resistance because all the electrons are collectively in the same state. In the ordinary flow of current you knock one electron or the other out of the regular flow, gradually deteriorating the general momentum. But here to get one electron away from what all the others are doing is very hard because of the tendency of all Bose particles to go in the same state. A current once started, just keeps on going forever.

It's also easy to understand that if you have a piece of metal in the superconducting state and turn on a magnetic field which isn't too strong (we won't go into the details of how strong), the magnetic field can't penetrate the metal. If, as you build up the magnetic field, any of it were to build up inside the metal, there would be a rate of change of flux which would produce an electric field, and an electric field would immediately generate a current which, by Lenz's law, would oppose the flux. Since all the electrons will move together, an infinitesimal electric field will generate enough current to oppose completely any applied magnetic field. So if you turn the field on after you've cooled a metal to the superconducting state, it will be excluded.

Even more interesting is a related phenomenon discovered experimentally by Meissner.⁸ If you have a piece of the metal at a high temperature (so that it is a normal conductor) and establish a magnetic field through it, and then you lower the temperature below the critical temperature (where the metal becomes a superconductor), *the field is expelled*. In other words, it starts up its own current—and in just the right amount to push the field out.

We can see the reason for that in the equations, and I'd like to explain how. Suppose that we take a piece of superconducting material which is in one lump. Then in a steady situation of any kind the divergence of the current must be zero because there's no place for it to go. It is convenient to choose to make the divergence of \mathbf{A} equal to zero. (I should explain why choosing this convention doesn't mean any loss of generality, but I don't want to take the time.) Taking the divergence of Eq. (21.18), then gives that the Laplacian of θ is equal to zero. One moment. What about the variation of ρ ? I forgot to mention an important point. There is a background of positive charge in this metal due to the atomic ions of the lattice. If the charge density ρ is uniform there is no net charge and no electric field. If there would be any accumulation of electrons in one region the charge wouldn't be neutralized and there would be a terrific repulsion pushing the electrons apart.[†] So in ordinary circumstances the charge density of the electrons in the superconductor is almost perfectly uniform—I can take ρ as a constant. Now the only way that $\nabla^2\theta$ can be zero everywhere inside the lump of metal is for θ to be a constant. And that means that there is no contribution to \mathbf{J} from p -momentum. Equation (21.18) then says that the current is proportional to ρ times \mathbf{A} . So everywhere in a lump of superconducting material the current is necessarily proportional to the vector potential:

$$\mathbf{J} = -\rho \frac{q}{m} \mathbf{A}. \quad (21.20)$$

Since ρ and q have the same (negative) sign, and since ρ is a constant, I can set $\rho q/m = -(\text{some constant})$; then

$$\mathbf{J} = -(\text{some constant})\mathbf{A}. \quad (21.21)$$

This equation was originally proposed by London and London⁹ to explain the experimental observations of superconductivity—long before the quantum mechanical origin of the effect was understood.

Now we can use Eq. (21.20) in the equations of electromagnetism to solve for the fields. The vector potential is related to the current density by

$$\nabla^2 \mathbf{A} = -\frac{1}{\epsilon_0 c^2} \mathbf{J}. \quad (21.22)$$

If I use Eq. (21.21) for \mathbf{J} , I have

$$\nabla^2 \mathbf{A} = \lambda^2 \mathbf{A}, \quad (21.23)$$

where λ^2 is just a new constant;

$$\lambda^2 = \rho \frac{q}{\epsilon_0 m c^2}. \quad (21.24)$$

We can now try to solve this equation for \mathbf{A} and see what happens in detail. For example, in one dimension Eq. (21.23) has exponential solutions of the form $e^{-\lambda x}$ and $e^{+\lambda x}$. These solutions mean that the vector potential must *decrease* exponentially as you go from the surface into the material. (It can't increase

⁸ W. Meissner and R. Ochsenfeld, *Naturwiss.* **21**, 787 (1933).

⁹ H. London and F. London, *Proc. Roy. Soc. (London)* **A149**, 71 (1935); *Physica* **2**, 341 (1935).

[†] Actually if the electric field were too strong, pairs would be broken up and the "normal" electrons created would move in to help neutralize any excess of positive charge. Still, it takes energy to make these normal electrons, so the main point is that a nearly uniform density ρ is highly favored energetically.

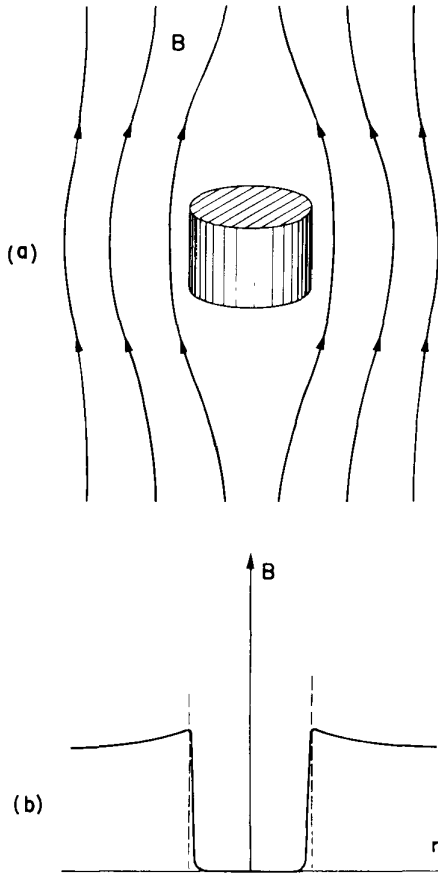


Fig. 21-3 (a) A superconducting cylinder is a magnetic field; (b) the magnetic field B as a function of r .

because there would be a blow up.) If the piece of metal is very large compared to $1/\lambda$, the field only penetrates to a thin layer at the surface—a layer about $1/\lambda$ in thickness. The entire remainder of the interior is free of field, as sketched in Fig. 21-3. This is the explanation of the Meissner effect.

How big is the distance λ ? Well, remember that r_0 , the “electromagnetic radius” of the electron (2.8×10^{-13} cm), is given by

$$mc^2 = \frac{q_e^2}{4\pi\epsilon_0 r_0}.$$

Also, remember that q in Eq. (21.24) is twice the charge on an electron, so

$$\frac{q}{\epsilon_0 mc^2} = \frac{8\pi r_0}{q_e}.$$

Writing ρ as $q_e N$, where N is the number of electrons per cubic centimeter, we have

$$\lambda^2 = 8\pi N r_0. \quad (21.25)$$

For a metal such as lead there are about 3×10^{22} atoms per cm^3 , so if each one contributed only one conduction electron, $1/\lambda$ would be about 2×10^{-5} cm. That gives you the order of magnitude

21-7 Flux quantization

The London equation (21.21) was proposed to account for the observed facts of superconductivity including the Meissner effect. In recent times, however, there have been some even more dramatic predictions. One prediction made by London was so peculiar that nobody paid much attention to it until recently. I will now discuss it. This time instead of taking a single lump, suppose we take a *ring* whose thickness is large compared to $1/\lambda$, and try to see what would happen if we started with a magnetic field through the ring, then cooled it to the superconducting state, and afterward removed the original source of \mathbf{B} . The sequence of events is sketched in Fig. 21-4. In the normal state there will be a field in the body of the ring as sketched in part (a) of the figure. When the ring is made superconducting, the field is forced outside of the *material* (as we have just seen). There will then be some flux through the hole of the ring as sketched in part (b). If the external field is now removed, the lines of field going through the hole are “trapped” as shown in part (c). The flux Φ through the center can’t decrease because $\partial\Phi/\partial t$ must be equal to the line integral of \mathbf{E} around the ring, which is zero in a superconductor. As the external field is removed a super current starts flowing around the ring to keep the flux through the ring a constant. (It’s the old eddy-current idea, only with zero resistance.) These currents will, however, all flow near the surface (down to a depth $1/\lambda$), as can be shown by the same kind of analysis that I made for the solid block. These currents can keep the magnetic field out of the body of the ring, and produce the permanently trapped magnetic field as well.

Now, however, there is an essential difference, and our equations predict a surprising effect. The argument I made above that θ must be a constant in a solid block *does not apply for a ring*, as you can see from the following arguments.

Well inside the body of the ring the current density \mathbf{J} is zero; so Eq. (21.18) gives

$$\hbar \nabla \theta = q \mathbf{A}. \quad (21.26)$$

Now consider what we get if we take the line integral of \mathbf{A} around a curve Γ , which goes around the ring near the center of its cross-section so that it never gets near the surface, as drawn in Fig. 21-5. From Eq. (21.26),

$$\hbar \oint \nabla \theta \cdot d\mathbf{s} = q \oint \mathbf{A} \cdot d\mathbf{s} \quad (21.27)$$

Now you know that the line integral of A around any loop is equal to the flux of B through the loop

$$\oint A \cdot ds = \Phi.$$

Equation (21.27) then becomes

$$\oint \nabla \theta \cdot ds = \frac{q}{\hbar} \Phi. \quad (21.28)$$

The line integral of a gradient from one point to another (say from point 1 to point 2) is the difference of the values of the function at the two points. Namely,

$$\int_1^2 \nabla \theta \cdot ds = \theta_2 - \theta_1.$$

If we let the two end points 1 and 2 come together to make a closed loop you might at first think that θ_2 would equal θ_1 , so that the integral in Eq. (21.28) would be zero. That would be true for a closed loop in a simply-connected piece of superconductor, but it is not necessarily true for a ring-shaped piece. The only physical requirement we can make is that *there can be only one value of the wave function for each point*. Whatever θ does as you go around the ring, when you get back to the starting point the θ you get must give the same value for the wave function

$$\psi = \sqrt{\rho} e^{i\theta}.$$

This will happen if θ changes by $2\pi n$, where n is any integer. So if we make one complete turn around the ring the left-hand side of Eq. (21.27) must be $\hbar \cdot 2\pi n$. Using Eq. (21.28), I get that

$$2\pi n \hbar = q \Phi. \quad (21.29)$$

The trapped flux must always be an integer times $2\pi\hbar/q$! If you would think of the ring as a classical object with an ideally perfect (that is, infinite) conductivity, you would think that whatever flux was initially found through it would just stay there—any amount of flux at all could be trapped. But the quantum-mechanical theory of superconductivity says that the flux can be zero, or $2\pi\hbar/q$, or $4\pi\hbar/q$, or $6\pi\hbar/q$, and so on, but no value in between. It must be a multiple of a basic quantum mechanical unit.

London¹⁰ predicted that the flux trapped by a superconducting ring would be quantized and said that the possible values of the flux would be given by Eq. (21.29) with q equal to the electronic charge. According to London the basic unit of flux should be $2\pi\hbar/q_e$, which is about 4×10^{-7} gauss \cdot cm². To visualize such a flux, think of a tiny cylinder a tenth of a millimeter in diameter; the magnetic field inside it when it contains this amount of flux is about one percent of the earth's magnetic field. It should be possible to observe such a flux by a sensitive magnetic measurement.

In 1961 such a quantized flux was looked for and found by Deaver and Fairbank¹¹ at Stanford University and at about the same time by Doll and Nabauer¹² in Germany.

In the experiment of Deaver and Fairbank, a tiny cylinder of superconductor was made by electroplating a thin layer of tin on a one-centimeter length of No. 56 (1.3×10^{-3} cm diameter) copper wire. The tin becomes superconducting below 3.8°K, while the copper remains a normal metal. The wire was put in a small controlled magnetic field, and the temperature reduced until the tin became superconducting. Then the external source of field was removed. You would

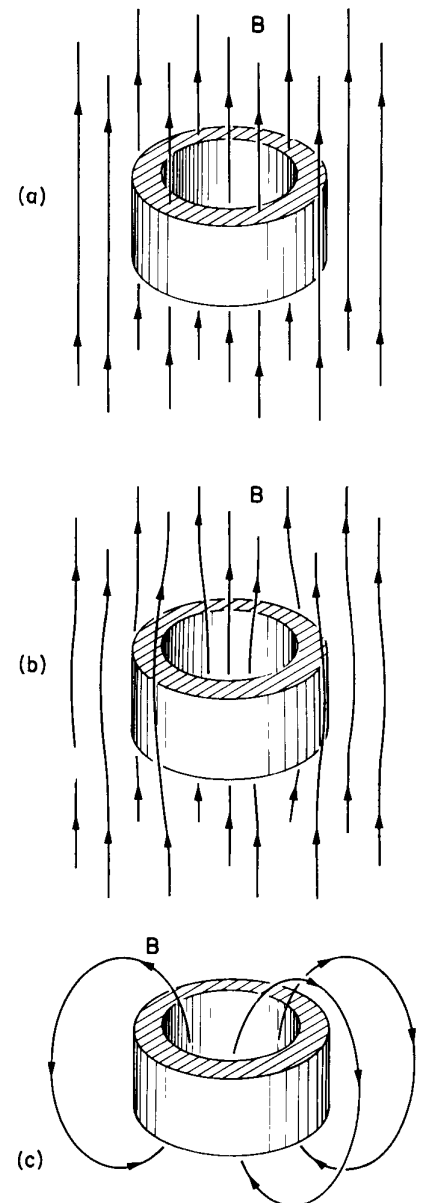


Fig. 21-4. A ring in a magnetic field: (a) in the normal state; (b) in the superconducting state; (c) after the external field is removed.

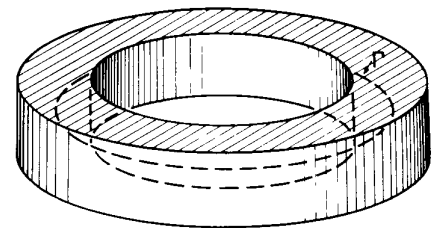


Fig. 21-5. The curve Γ inside a superconducting ring

¹⁰ F. London, *Superfluids*; John Wiley and Sons, Inc., New York, 1950, Vol. I, p. 152

¹¹ B. S. Deaver, Jr., and W. M. Fairbank, *Phys. Rev. Letters* 7, 43 (1961).

¹² R. Doll and M. Nabauer, *Phys. Rev. Letters* 7, 51 (1961).

expect this to generate a current by Lenz's law so that the flux inside would not change. The little cylinder should now have magnetic moment proportional to the flux inside. The magnetic moment was measured by jiggling the wire up and down (like the needle on a sewing machine, but at the rate of 100 cycles per second) inside a pair of little coils at the ends of the tin cylinder. The induced voltage in the coils was then a measure of the magnetic moment.

When the experiment was done by Deaver and Fairbank, they found that the flux was quantized, *but that the basic unit was only one-half as large as London had predicted.* Doll and Nabauer got the same result. At first this was quite mysterious,[†] but we now understand why it should be so. According to the Bardeen, Cooper, and Schrieffer theory of superconductivity, the q which appears in Eq. (21.29) is the charge of a *pair* of electrons and so is equal to $2q_e$. The basic flux unit is

$$\Phi_0 = \frac{\pi\hbar}{q_e} \approx 2 \times 10^{-7} \text{ gauss-cm} \quad (21.30)$$

or one-half the amount predicted by London. Everything now fits together, and the measurements show the existence of the predicted purely quantum-mechanical effect on a large scale.

21-8 The dynamics of superconductivity

The Meissner effect and the flux quantization are two confirmations of our general ideas. Just for the sake of completeness I would like to show you what the complete equations of a superconducting fluid would be from this point of view—it is rather interesting. Up to this point I have only put the expression for ψ into equations for charge density and current. If I put it into the complete Schrödinger equation I get equations for ρ and θ . It should be interesting to see what develops, because here we have a “fluid” of electron pairs with a charge density ρ and a mysterious θ —we can try to see what kind of equations we get for such a “fluid”! So we substitute the wave function of Eq. (21.17) into the Schrödinger equation (21.3) and remember that ρ and θ are real functions of x , y , and z . If we separate real and imaginary parts we obtain then two equations. To write them in a shorter form I will—following Eq. (21.19)—write

$$\frac{\hbar}{m} \nabla \theta - \frac{q}{m} \mathbf{A} = \mathbf{v}. \quad (21.31)$$

One of the equations I get is then

$$\frac{\partial \rho}{\partial t} = \nabla \cdot \rho \mathbf{v}. \quad (21.32)$$

Since $\rho \mathbf{v}$ is first \mathbf{J} , this is just the continuity equation once more. The other equation I obtain tells how θ varies; it is

$$\hbar \frac{\partial \theta}{\partial t} = -\frac{m}{2} v^2 + q\phi - \frac{\hbar^2}{2m} \left\{ \frac{1}{\sqrt{\rho}} \nabla^2 (\sqrt{\rho}) \right\}. \quad (21.33)$$

Those who are thoroughly familiar with hydrodynamics (of which I'm sure few of you are) will recognize this as the equation of motion for an electrically charged fluid if we identify $\hbar\theta$ as the “velocity potential”—except that the last term, which should be the energy of compression of the fluid, has a rather strange dependence on the density ρ . In any case, the equation says that rate of change of the quantity $\hbar\theta$ is given by a kinetic energy term, $\frac{1}{2}mv^2$, plus a potential energy term, $q\phi$, with an additional term, containing the factor \hbar^2 , which we could call a “quantum mechanical energy.” We have seen that inside a superconductor ρ is kept very

[†] It has once been suggested by Onsager that this might happen (see F. London, Ref. 10), although no one else ever understood why.

uniform by the electrostatic forces, so this term can almost certainly be neglected in every practical application provided we have only one superconducting region. If we have a boundary between two superconductors (or other circumstances in which the value of ρ may change rapidly) this term can become important.

For those who are not so familiar with the equations of hydrodynamics, I can rewrite Eq. (21.33) in a form that makes the physics more apparent by using Eq. (21.31) to express θ in terms of \mathbf{v} . Taking the gradient of the whole of Eq. (21.33) and expressing $\nabla\theta$ in terms of \mathbf{A} and \mathbf{v} by using (21.31), I get

$$\frac{\partial \mathbf{v}}{\partial t} = \frac{q}{m} \left(-\nabla\phi - \frac{\partial \mathbf{A}}{\partial t} \right) - \mathbf{v} \times (\nabla \times \mathbf{v}) - (\mathbf{v} \times \nabla)\mathbf{v} - \nabla \frac{\hbar^2}{2m} \left(\frac{1}{\sqrt{\rho}} \nabla^2 \sqrt{\rho} \right). \quad (21.34)$$

What does this equation mean? First, remember that

$$-\nabla\phi - \frac{\partial \mathbf{A}}{\partial t} = \mathbf{E}. \quad (21.35)$$

Next, notice that if I take the curl of Eq. (21.19), I get

$$\nabla \times \mathbf{v} = -\frac{q}{m} \nabla \times \mathbf{A}, \quad (21.36)$$

since the curl of a gradient is always zero. But $\nabla \times \mathbf{A}$ is the magnetic field \mathbf{B} , so the first two terms can be written as

$$\frac{q}{m} (\mathbf{E} + \mathbf{v} \times \mathbf{B}).$$

Finally, you should understand that $\partial \mathbf{v} / \partial t$ stands for the rate of change of the velocity of the fluid at a point. If you concentrate on a particular particle, its acceleration is the *total* derivative of \mathbf{v} (or, as it is sometimes called in fluid dynamics, the “comoving acceleration”), which is related to $\partial \mathbf{v} / \partial t$ by¹³

$$\left. \frac{d\mathbf{v}}{dt} \right|_{\text{comoving}} = \frac{\partial \mathbf{v}}{\partial t} + (\mathbf{v} \cdot \nabla)\mathbf{v}. \quad (21.37)$$

This extra term also appears as the third term on the right side of Eq. (21.25). Taking it to the left side, I can write Eq. (21.25) in the following way:

$$m \left. \frac{d\mathbf{v}}{dt} \right|_{\text{comoving}} = q(\mathbf{E} + \mathbf{v} \times \mathbf{B}) - \nabla \frac{\hbar^2}{2} \left(\frac{1}{\sqrt{\rho}} \nabla^2 \sqrt{\rho} \right). \quad (21.38)$$

We also have from Eq. (21.36) that

$$\nabla \times \mathbf{v} = -\frac{q}{m} \mathbf{B}. \quad (21.39)$$

These two equations are the equations of motion of the superconducting electron fluid. The first equation is just Newton’s law for a charged fluid in an electromagnetic field. It says that the acceleration of each particle of the fluid whose charge is q comes from the ordinary Lorentz force $q(\mathbf{E} + \mathbf{v} \times \mathbf{B})$ plus an additional force, which is the gradient of some mystical quantum mechanical potential—a force which is not very big except at the junction between two superconductors. The second equation says that the fluid is “ideal”—the curl of \mathbf{v} has zero divergence (the divergence of \mathbf{B} is always zero). That means that the velocity can be expressed in terms of velocity potential. Ordinarily one writes that $\nabla \times \mathbf{v} = 0$ for an ideal fluid, but for an *ideal charged fluid in a magnetic field*, this gets modified to Eq. (21.40).

So, Schrödinger’s equation for the electron pairs in a superconductor gives us the equations of motion of an electrically charged ideal fluid. Superconductivity is the same as the problem of the hydrodynamics of a charged liquid. If you want

¹³ See Volume II, Section 40–2.

to solve any problem about superconductors you take these equations for the fluid [or the equivalent pair, Eqs. (21.32) and (21.33)], and combine them with Maxwell's equations to get the fields. (The charges and currents you use to get the fields must, of course, include the ones from the superconductor as well as from the external sources.)

Incidentally, I believe that Eq. (21.38) is not quite correct, but ought to have an additional term involving the density. This new term does not depend on quantum mechanics, but comes from the ordinary energy associated with variations of density. Just as in an ordinary fluid there should be a potential energy density proportional to the square of the deviation of ρ from ρ_0 , the undisturbed density (which is, here, also equal to the charge density of the crystal lattice). Since there will be forces proportional to the gradient of this energy, there should be another term in Eq. (21.38) of the form: $(\text{const}) \nabla(\rho - \rho_0)^2$. This term did not appear from the analysis because it comes from the interactions between particles, which I neglected in using an independent-particle approximation. It is, however, just the force I referred to when I made the qualitative statement that electrostatic forces would tend to keep ρ nearly constant inside a superconductor.

21-9 The Josephson junction

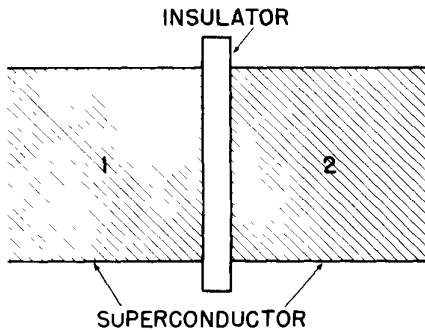


Fig. 21-6. Two superconductors separated by a thin insulator.

I would like to discuss next a very interesting situation that was noticed by Josephson¹⁴ while analyzing what might happen at a junction between two superconductors. Suppose we have two superconductors which are connected by a thin layer of insulating material as in Fig. 21-6. Such an arrangement is now called a “Josephson junction.” If the insulating layer is thick, the electrons can’t get through; but if the layer is thin enough, there can be an appreciable quantum mechanical amplitude for electrons to jump across. This is just another example of the quantum-mechanical penetration of a barrier. Josephson analyzed this situation and discovered that a number of strange phenomenon should occur.

In order to analyze such a junction I’ll call the amplitude to find an electron on one side, ψ_1 , and the amplitude to find it on the other, ψ_2 . In the superconducting state the wave function, ψ_1 is the common wave function of all the electrons on one side, and ψ_2 is the corresponding function on the other side. I could do this problem for different kinds of superconductors, but let us take a very simple situation in which the material is the same on both sides so that the junction is symmetrical and simple. Also, for a moment let there be no magnetic field. Then the two amplitudes should be related in the following way:

$$i\hbar \frac{\partial \psi_1}{\partial t} = U_1 \psi_1 + K \psi_2,$$

$$i\hbar \frac{\partial \psi_2}{\partial t} = U_2 \psi_2 + K \psi_1.$$

The constant K is a characteristic of the junction. If K were zero, these two equations would just describe the lowest energy state—with energy U —of each superconductor. But there is coupling between the two sides by the amplitude K that there may be leakage from one side to the other. (It is just the “flip-flop” amplitude of a two-state system.) If the two sides are identical, U_1 would equal U_2 and I could just subtract them off. But now suppose that we connect the two superconducting regions to the two terminals of a battery so that there is a potential difference V across the junction. Then $U_1 - U_2 = qV$. I can, for convenience, define the zero of energy to be halfway between, then the two equations are

$$i\hbar \frac{\partial \psi_1}{\partial t} = \frac{qV}{2} \psi_1 + K \psi_2,$$

$$i\hbar \frac{\partial \psi_2}{\partial t} = -\frac{qV}{2} \psi_2 + K \psi_1. \quad (21.40)$$

¹⁴ B. D. Josephson, *Physics Letters* 1, 251 (1962).

These are the standard equations for two quantum mechanical states coupled together. This time, let's analyze these equations in another way. Let's make the substitutions

$$\begin{aligned}\psi_1 &= \sqrt{\rho_1} e^{i\theta_1}, \\ \psi_2 &= \sqrt{\rho_2} e^{i\theta_2},\end{aligned}\tag{21.41}$$

where θ_1 and θ_2 are the phases on the two sides of the junction and ρ_1 and ρ_2 are the density of electrons at those two points. Remember that in actual practice ρ_1 and ρ_2 are almost exactly the same and are equal to ρ_0 , the normal density of electrons in the superconducting material. Now if you substitute these equations for ψ_1 and ψ_2 into (21.40), you get four equations by equating the real and imaginary parts in each case. Letting $(\theta_2 - \theta_1) = \delta$, for short, the result is

$$\dot{\rho}_1 = + \frac{2}{\hbar} K \sqrt{\rho_2 \rho_1} \sin \delta,\tag{21.42}$$

$$\dot{\rho}_2 = - \frac{2}{\hbar} K \sqrt{\rho_2 \rho_1} \sin \delta,$$

$$\dot{\theta}_1 = + \frac{K}{\hbar} \sqrt{\frac{\rho_2}{\rho_1}} \cos \delta - \frac{qV}{2\hbar},\tag{21.43}$$

$$\dot{\theta}_2 = + \frac{K}{\hbar} \sqrt{\frac{\rho_1}{\rho_2}} \cos \delta + \frac{qV}{2\hbar}.$$

The first two equations say that $\dot{\rho}_1 = -\dot{\rho}_2$. "But," you say, "they must both be zero if ρ_1 and ρ_2 are both constant and equal to zero." Not quite. These equations are not the whole story. They say what $\dot{\rho}_1$ and $\dot{\rho}_2$ would be *if there were no extra electric forces* due to an unbalance between the electron fluid and the background of positive ions. They tell how the densities would *start* to change, and therefore describe the kind of current that would begin to flow. This current from side 1 to side 2 would be just $\dot{\rho}_1$ (or $-\dot{\rho}_2$), or

$$J = \frac{2K}{\hbar} \sqrt{\rho_1 \rho_2} \sin \delta.\tag{21.44}$$

Such a current would soon charge up side 2, *except* that we have forgotten that the two sides are connected by wires to the battery. The current that flows will not charge up region 2 (or discharge region 1) because currents will flow to keep the potential constant. These currents from the battery have not been included in our equations. When they are included, ρ_1 and ρ_2 do not in fact change, but the current across the junction is still given by Eq. (21.44).

Since ρ_1 and ρ_2 do remain constant and equal to ρ_0 , let's set $2K\rho_0/\hbar = J_0$, and write

$$J = J_0 \sin \delta.\tag{21.45}$$

J_0 , like K , is then a number which is a characteristic of the particular junction.

The other pair of equations (21.43) tells us about θ_1 and θ_2 . We are interested in the difference $\delta = \theta_2 - \theta_1$ to use Eq. (21.45); what we get is

$$\dot{\delta} = \dot{\theta}_2 - \dot{\theta}_1 = \frac{qV}{\hbar}.\tag{21.46}$$

That means that we can write

$$\delta(t) = \delta_0 + \frac{q}{\hbar} \int V(t) dt,\tag{21.47}$$

where δ_0 is the value of δ at $t = 0$. Remember also that q is the charge of a pair, namely, $q = 2q_e$. In Eqs. (21.45) and (21.47) we have an important result, the general theory of the Josephson junction.

Now what are the consequences? First, put on a dc voltage. If you put on a dc voltage, V_0 , the argument of the sine becomes $(\delta_0 + (q/\hbar)V_0 t)$. Since \hbar is a small number (compared to ordinary voltage and times), the sine oscillates rather rapidly and the net current is nothing. (In practice, since the temperature is not zero, you would get a small current due to the conduction by “normal” electrons.) On the other hand if you have *zero* voltage across the junction, you can get a current! With no voltage the current can be any amount between $+J_0$ and $-J_0$ (depending on the value of δ_0). But try to put a voltage across it and the current goes to zero. This strange behavior has recently been observed experimentally.¹⁵

There is another way of getting a current—by applying a voltage at a very high frequency in addition to a dc voltage. Let

$$V = V_0 + v \cos \omega t,$$

where $v \ll V$. Then $\delta(t)$ is

$$\delta_0 + \frac{q}{\hbar} V_0 t + \frac{q}{\hbar} \frac{v}{\omega} \sin \omega t.$$

Now for Δx small,

$$\sin(x + \Delta x) \approx \sin x + \Delta x \cos x.$$

Using this approximation for $\sin \delta$, I get

$$J = J_0 \left[\sin \left(\delta_0 + \frac{q}{\hbar} V_0 t \right) + \frac{q}{\hbar} \frac{v}{\omega} \sin \omega t \cos \left(\delta_0 + \frac{q}{\hbar} V_0 t \right) \right].$$

The first term is zero on the average, but the second term is not if

$$\omega = \frac{q}{\hbar} V_0.$$

There should be a current if the ac voltage has just this frequency. Shapiro¹⁶ claims to have observed such a resonance effect.

If you look up papers on the subject you will find that they often write the formula for the current as

$$J = J_0 \sin \left(\delta_0 + \frac{2q_e}{\hbar} \int A \cdot ds \right), \quad (21.48)$$

where the integral is to be taken across the junction. The reason for this is that when there's a vector potential across the junction the flip-flop amplitude is modified in phase in the way that we explained earlier. If you chase that extra phase through, it comes out as given above.

Finally, I would like to describe a very dramatic and interesting experiment which has recently been made on the interference of the currents from each of two junctions. In quantum mechanics we're used to the interference between amplitudes from two different slits. Now we're going to do the interference between two junctions caused by the difference in the phase of the arrival of the currents through two different paths. In Fig. 21-7, I show two different junctions, “a” and “b”, connected in parallel. The ends, P and Q , are connected to our electrical instruments which measure any current flow. The external current, J_{total} , will be the sum of the currents through the two junctions. Let J_a and J_b be the currents through the two junctions, and let their phases be δ_a and δ_b . Now the phase difference of the wave functions between P and Q must be the same whether you go on one route or the other. Along the route through junction “a”, the phase difference between P and Q is δ_a plus the line integral of the vector potential along the upper route:

$$\Delta \text{Phase}_{P \rightarrow Q} = \delta_a + \frac{2q_e}{\hbar} \int_{\text{upper}} A \cdot ds. \quad (21.49)$$

¹⁵ P. W. Anderson and J. M. Rowell, *Phys. Rev. Letters* **10**, 230 (1963).

¹⁶ S. Shapiro, *Phys. Rev. Letters* **11**, 80 (1963).

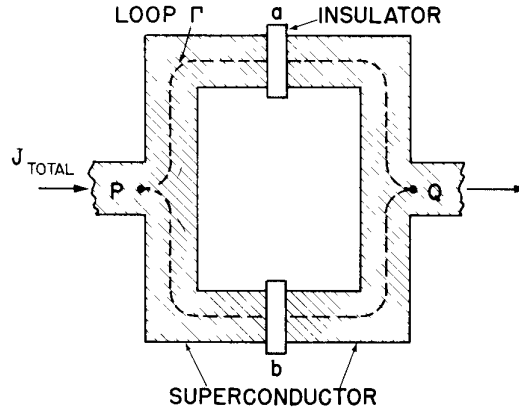


Fig. 21-7. Two Josephson junctions in parallel.

Why? Because the phase θ is related to A by Eq. (21.26). If you integrate that equation along some path, the left-hand side gives the phase change, which is then just proportional to the line integral of A , as we have written here. The phase change along the lower route can be written similarly

$$\Delta\text{Phase}_{P \rightarrow Q} = \delta_b + \frac{2q_e}{\hbar} \int_{\text{lower}} A \cdot ds. \quad (21.50)$$

These two must be equal; and if I subtract them I get that the difference of the deltas must be the line integral of A around the circuit:

$$\delta_b - \delta_a = \frac{2q_e}{\hbar} \oint_{\Gamma} A \cdot ds.$$

Here the integral is around the closed loop Γ of Fig. 21-7 which circles through both junctions. The integral over A is the magnetic flux Φ through the loop. So the two δ 's are going to differ by $2q_e/\hbar$ times the magnetic flux Φ which passes between the two branches of the circuit:

$$\delta_b - \delta_a = \frac{2q_e}{\hbar} \Phi. \quad (21.51)$$

I can control this phase difference by changing the magnetic field on the circuit, so I can adjust the differences in phases and see whether or not the total current that flows through the two junctions shows any interference of the two parts. The total current will be the sum of J_a and J_b . For convenience, I will write

$$\delta_a = \delta_0 + \frac{q_e}{\hbar} \Phi, \quad \delta_b = \delta_0 - \frac{q_e}{\hbar} \Phi.$$

Then,

$$\begin{aligned} J_{\text{total}} &= J_0 \left\{ \sin \left(\delta_0 + \frac{q_e}{\hbar} \Phi \right) + \sin \left(\delta_0 - \frac{q_e}{\hbar} \Phi \right) \right\} \\ &= J_0 \sin \delta_0 \cos \frac{q_e \Phi}{\hbar}. \end{aligned} \quad (21.52)$$

Now we don't know anything about δ_0 , and nature can adjust that anyway she wants depending on the circumstances. In particular, it will depend on the external voltage we apply to the junction. No matter what we do, however, $\sin \delta_0$ can never get bigger than 1. So the *maximum* current for any given Φ is given by

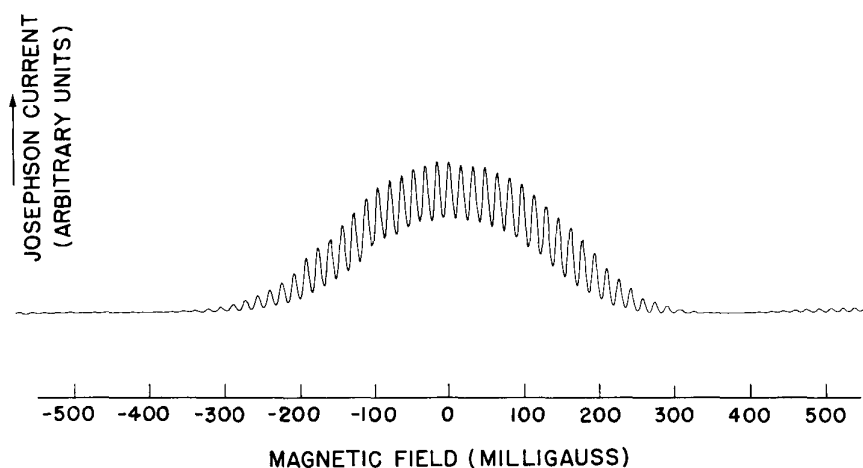
$$J_{\text{max}} = J_0 \left| \cos \frac{q_e \Phi}{\hbar} \right|.$$

This maximum current will vary with Φ and will itself have maxima whenever

$$\Phi = n \frac{\pi \hbar}{q_e},$$

with n some integer. That is to say that the current takes on its maximum values where the flux linkage has just those quantized values we found in Eq. (21.30)

Fig. 21-8. A recording of the current through a pair of Josephson junctions as a function of the magnetic field in the region between the two junctions (see Fig. 21-7). [This recording was provided by R. C. Jaklevic, J. Lambe, A. H. Silver, and J. E. Mercereau of the Scientific Laboratory, Ford Motor Company]



The Josephson current through a double junction was recently measured¹⁷ as a function of the magnetic field in the area between the junctions. The results are shown in Fig. 21-8. There is a general background of current from various effects we have neglected, but the rapid oscillations of the current with changes in the magnetic field are due to the interference term $\cos q_e \Phi / h$ of Eq. (21.52).

One of the intriguing questions about quantum mechanics is the question of whether the vector potential exists in a place where there's no field.¹⁸ This experiment I have just described has also been done with a tiny solenoid between the two junctions so that the only significant magnetic \mathbf{B} field is inside the solenoid and a negligible amount is on the superconducting wires themselves. Yet it is reported that the amount of current depends oscillatorily on the flux of magnetic field inside that solenoid even though that field never touches the wires—another demonstration of the “physical reality” of the vector potential.¹⁹

I don't know what will come next. But look what can be done. First, notice that the interference between two junctions can be used to make a sensitive magnetometer. If a pair of junctions is made with an enclosed area of, say, 1 mm^2 , the maxima in the curve of Fig. 21-8 would be separated by 2×10^{-6} gauss. It is certainly possible to tell when you are 1/10 of the way between two peaks; so it should be possible to use such a junction to measure magnetic fields as small as 2×10^{-7} gauss—or to measure larger fields to such a precision. One should be able to go even farther. Suppose for example we put a set of 10 or 20 junctions close together and equally spaced. Then we can have the interference between 10 or 20 slits and as we change the magnetic field we will get very sharp maxima and minima. Instead of a 2-slit interference we can have a 20- or perhaps even a 100-slit interferometer for measuring the magnetic field. Perhaps we can predict that the measurement of magnetic fields will—by using the effects of quantum-mechanical interference—eventually become almost as precise as the measurement of wavelength of light.

These then are some illustrations of things that are happening in modern times—the transistor, the laser, and now these junctions, whose ultimate practical applications are still not known. The quantum mechanics which was discovered in 1926 has had nearly 40 years of development, and rather suddenly it has begun to be exploited in many practical and real ways. We are really getting control of nature on a very delicate and beautiful level.

I am sorry to say, gentlemen, that to participate in this adventure it is absolutely imperative that you learn quantum mechanics as soon as possible. It was our hope that in this course we would find a way to make comprehensible to you at the earliest possible moment the mysteries of this part of physics

¹⁷ Jaklevic, Lambe, Silver, and Mercereau, *Phys. Rev. Letters* **12**, 159 (1964).

¹⁸ Jaklevic, Lambe, Silva, and Mercereau, *Phys. Rev. Letters* **12**, 274 (1964).

¹⁹ See Volume II, Chapter 15, Section 15-5.

Feynman's Epilogue

Well, I've been talking to you for two years and now I'm going to quit. In some ways I would like to apologize, and other ways not. I hope—in fact, I know—that two or three dozen of you have been able to follow everything with great excitement, and have had a good time with it. But I also know that “the powers of instruction are of very little efficacy except in those happy circumstances in which they are practically superfluous.” So, for the two or three dozen who have understood everything, may I say I have done nothing but shown you the things. For the others, if I have made you hate the subject, I'm sorry. I never taught elementary physics before, and I apologize. I just hope that I haven't caused a serious trouble to you, and that you do not leave this exciting business. I hope that someone else can teach it to you in a way that doesn't give you indigestion, and that you will find someday that, after all, it isn't as horrible as it looks.

Finally, may I add that the main purpose of my teaching has not been to prepare you for some examination—it was not even to prepare you to serve industry or the military. I wanted most to give you some appreciation of the wonderful world and the physicist's way of looking at it, which, I believe, is a major part of the true culture of modern times. (There are probably professors of other subjects who would object, but I believe that they are completely wrong.)

Perhaps you will not only have some appreciation of this culture; it is even possible that you may want to join in the greatest adventure that the human mind has ever begun.