

Simple Nature

Crowell



Simple Nature

An Introduction to Physics for Engineering
and Physical Science Students

Benjamin Crowell

www.lightandmatter.com



Fullerton, California
www.lightandmatter.com

Copyright ©2001-2008 Benjamin Crowell

rev. May 16, 2020

Permission is granted to copy, distribute and/or modify this document under the terms of the Creative Commons Attribution Share-Alike License, which can be found at creativecommons.org. The license applies to the entire text of this book, plus all the illustrations that are by Benjamin Crowell. (At your option, you may also copy this book under the GNU Free Documentation License version 1.2, with no invariant sections, no front-cover texts, and no back-cover texts.) All the illustrations are by Benjamin Crowell except as noted in the photo credits or in parentheses in the caption of the figure. This book can be downloaded free of charge from www.lightandmatter.com in a variety of formats, including editable formats.

Brief Contents

0	Introduction and Review	13
1	Conservation of Mass	55
2	Conservation of Energy	73
3	Conservation of Momentum	131
4	Conservation of Angular Momentum	251
5	Thermodynamics	307
6	Waves	353
7	Relativity	397
8	Atoms and Electromagnetism	473
9	Circuits	531
10	Fields	579
11	Electromagnetism	675
12	Optics	765
13	Quantum Physics	857
14	Additional Topics in Quantum Physics	959

Contents

0 Introduction and Review

0.1 Introduction and review	13
The scientific method, 13.—What is physics?, 16.—How to learn physics, 19.—Velocity and acceleration, 21.—Self-evaluation, 23.—Basics of the metric system, 24.—Less common metric prefixes, 27.—Scientific notation, 27.—Conversions, 28.—Significant figures, 30.—A note about diagrams, 32.	
0.2 Scaling and order-of-magnitude estimates	34
Introduction, 34.—Scaling of area and volume, 35.—Order-of-magnitude estimates, 43.	
Problems	47

1 Conservation of Mass

1.1 Mass	55
Problem-solving techniques, 58.—Delta notation, 59.	
1.2 Equivalence of gravitational and inertial mass	60
1.3 Galilean relativity	62
Applications of calculus, 67.	
1.4 A preview of some modern physics	68
Problems	70

2 Conservation of Energy

2.1 Energy	73
The energy concept, 73.—Logical issues, 75.—Kinetic energy, 76.—Power, 80.—Gravitational energy, 81.—Equilibrium and stability, 86.—Predicting the direction of motion, 89.	
2.2 Numerical techniques	91
2.3 Gravitational phenomena	96
Kepler's laws, 96.—Circular orbits, 98.—The sun's gravitational field, 99.—Gravitational energy in general, 99.—The shell theorem, 102.—★Evidence for repulsive gravity, 108.	
2.4 Atomic phenomena	109
Heat is kinetic energy., 110.—All energy comes from particles moving or interacting., 111.—Applications, 113.	
2.5 Oscillations	115
Problems	120
Exercises	128

3 Conservation of Momentum

3.1 Momentum in one dimension	132
Mechanical momentum, 132.—Nonmechanical momentum, 135.—	

Momentum compared to kinetic energy, 136.—Collisions in one dimension, 138.—The center of mass, 142.—The center of mass frame of reference, 147.—Totally inelastic collisions, 148.	
3.2 Force in one dimension	149
Momentum transfer, 149.—Newton’s laws, 150.—What force is not, 153.—Forces between solids, 155.—Fluid friction, 159.—Analysis of forces, 160.—Transmission of forces by low-mass objects, 162.—Work, 164.—Simple Machines, 171.—Force related to interaction energy, 172.	
3.3 Resonance	175
Damped, free motion, 176.—The quality factor, 179.—Driven motion, 180.	
3.4 Motion in three dimensions	191
The Cartesian perspective, 191.—Rotational invariance, 195.—Vectors, 197.—Calculus with vectors, 212.—The dot product, 216.—Gradients and line integrals (optional), 219.	
Problems	222
Exercises	244

4 Conservation of Angular Momentum

4.1 Angular momentum in two dimensions	251
Angular momentum, 251.—Application to planetary motion, 256.—Two theorems about angular momentum, 257.—Torque, 260.—Applications to statics, 264.—Proof of Kepler’s elliptical orbit law, 268.	
4.2 Rigid-body rotation	271
Kinematics, 271.—Relations between angular quantities and motion of a point, 272.—Dynamics, 274.—Iterated integrals, 276.—Finding moments of inertia by integration, 279.	
4.3 Angular momentum in three dimensions	284
Rigid-body kinematics in three dimensions, 284.—Angular momentum in three dimensions, 286.—Rigid-body dynamics in three dimensions, 291.	
Problems	294
Exercises	305

5 Thermodynamics

5.1 Pressure, temperature, and heat	308
Pressure, 308.—Temperature, 312.—Heat, 315.	
5.2 Microscopic description of an ideal gas	316
Evidence for the kinetic theory, 316.—Pressure, volume, and temperature, 317.	
5.3 Entropy as a macroscopic quantity	320
Efficiency and grades of energy, 320.—Heat engines, 321.—Entropy, 322.	
5.4 Entropy as a microscopic quantity	326
A microscopic view of entropy, 326.—Phase space, 328.—Microscopic definitions of entropy and temperature, 329.—Equipartition, 333.—The arrow of time, or “this way to the Big Bang”, 337.—Quantum	

mechanics and zero entropy, 339.—Summary of the laws of thermodynamics, 339.	
5.5 More about heat engines	340
Problems	347

6 Waves

6.1 Free waves	354
Wave motion, 354.—Waves on a string, 360.—Sound and light waves, 363.—Periodic waves, 365.—The Doppler effect, 368.	
6.2 Bounded waves	374
Reflection, transmission, and absorption, 374.—Quantitative treatment of reflection, 379.—Interference effects, 382.—Waves bounded on both sides, 384.—★Some technical aspects of reflection, 389.	
Problems	392

7 Relativity

7.1 Time is not absolute	397
The correspondence principle, 397.—Causality, 397.—Time distortion arising from motion and gravity, 398.	
7.2 Distortion of space and time	400
The Lorentz transformation, 400.—The γ factor, 405.—The universal speed c , 411.—No action at a distance, 416.—The light cone, 419.—★The spacetime interval, 420.—★Four-vectors and the inner product, 425.—★Doppler shifts of light and addition of velocities, 426.	
7.3 Dynamics	429
Momentum, 429.—Equivalence of mass and energy, 433.—★The energy-momentum four-vector, 437.—★Proofs, 440.	
7.4 ★General relativity	443
Our universe isn't Euclidean, 443.—The equivalence principle, 446.—Black holes, 450.—Cosmology, 453.	
Problems	457
Exercises	465

8 Atoms and Electromagnetism

8.1 The electric glue	473
The quest for the atomic force, 474.—Charge, electricity and magnetism, 475.—Atoms, 480.—Quantization of charge, 485.—The electron, 488.—The raisin cookie model of the atom, 492.	
8.2 The nucleus	494
Radioactivity, 494.—The planetary model, 497.—Atomic number, 501.—The structure of nuclei, 506.—The strong nuclear force, alpha decay and fission, 509.—The weak nuclear force; beta decay, 512.—Fusion, 514.—Nuclear energy and binding energies, 517.—Biological effects of ionizing radiation, 518.—★The creation of the elements, 523.	
Problems	525
Exercises	529

9 Circuits

9.1 Current and voltage	532
Current, 532.—Circuits, 535.—Voltage, 536.—Resistance, 541.— Current-conducting properties of materials, 550.	
9.2 Parallel and series circuits	554
Schematics, 554.—Parallel resistances and the junction rule, 555.— Series resistances, 559.	
Problems	566
Exercises	574

10 Fields

10.1 Fields of force	579
Why fields?, 579.—The gravitational field, 581.—The electric field, 585.	
10.2 Potential related to field	592
One dimension, 592.—Two or three dimensions, 595.	
10.3 Fields by superposition	597
Electric field of a continuous charge distribution, 597.—The field near a charged surface, 603.	
10.4 Energy in fields	606
Electric field energy, 606.—Gravitational field energy, 611.—Magnetic field energy, 611.	
10.5 LRC circuits	613
Capacitance and inductance, 613.—Oscillations, 617.—Voltage and current, 619.—Decay, 624.—Review of complex numbers, 627.— Euler's formula, 629.—Impedance, 630.—Power, 634.—Impedance matching, 637.—Impedances in series and parallel, 639.	
10.6 Fields by Gauss' law	641
Gauss' law, 641.—Additivity of flux, 645.—Zero flux from outside charges, 645.—Proof of Gauss' theorem, 649.—Gauss' law as a fundamental law of physics, 649.—Applications, 650.	
10.7 Gauss' law in differential form	653
Gauss's law as a local law, 653.—Poisson's equation and Laplace's equation, 657.—The method of images, 657.	
Problems	659
Exercises	671

11 Electromagnetism

11.1 More about the magnetic field	675
Magnetic forces, 675.—The magnetic field, 679.—Some applica- tions, 683.—No magnetic monopoles, 685.—Symmetry and hand- edness, 687.	
11.2 Magnetic fields by superposition	689
Superposition of straight wires, 689.—Energy in the magnetic field, 693.—Superposition of dipoles, 693.—The g factor (optional), 697.— The Biot-Savart law (optional), 698.	
11.3 Magnetic fields by Ampère's law	702
Ampère's law, 702.—A quick and dirty proof, 704.—Maxwell's equations for static fields, 705.	

11.4 Ampère's law in differential form (optional)	707
The curl operator, 707.—Properties of the curl operator, 708.	
11.5 Induced electric fields	713
Faraday's experiment, 713.—Why induction?, 717.—Faraday's law, 719.	
11.6 Maxwell's equations	724
Induced magnetic fields, 724.—Light waves, 726.	
11.7 Electromagnetic properties of materials	736
Conductors, 736.—Dielectrics, 737.—Magnetic materials, 739.	
Problems	747
Exercises	762

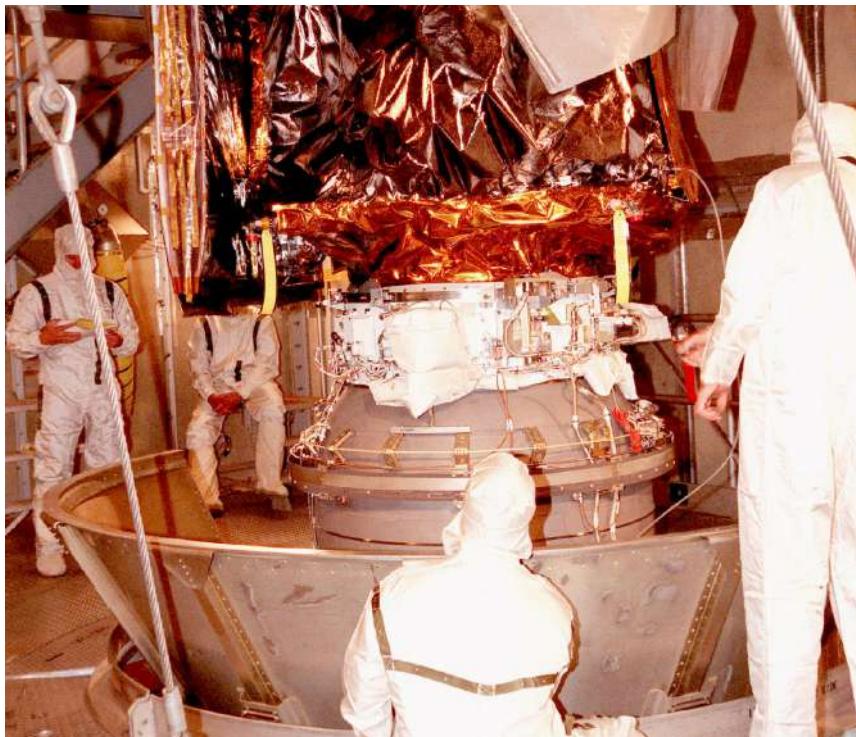
12 Optics

12.1 The ray model of light	765
The nature of light, 766.—Interaction of light with matter, 769.—The ray model of light, 771.—Geometry of specular reflection, 774.—★The principle of least time for reflection, 778.	
12.2 Images by reflection	780
A virtual image, 780.—Curved mirrors, 783.—A real image, 784.—Images of images, 785.	
12.3 Images, quantitatively	790
A real image formed by a converging mirror, 790.—Other cases with curved mirrors, 794.—★Aberrations, 798.	
12.4 Refraction	801
Refraction, 802.—Lenses, 808.—★The lensmaker's equation, 810.—Dispersion, 811.—★The principle of least time for refraction, 811.—★Microscopic description of refraction, 812.	
12.5 Wave optics	814
Diffraction, 814.—Scaling of diffraction, 816.—The correspondence principle, 816.—Huygens' principle, 817.—Double-slit diffraction, 818.—Repetition, 822.—Single-slit diffraction, 823.—Coherence, 825.—★The principle of least time, 827.	
Problems	829
Exercises	847

13 Quantum Physics

13.1 Rules of randomness	857
Randomness isn't random., 858.—Calculating randomness, 859.—Probability distributions, 863.—Exponential decay and half-life, 865.—Applications of calculus, 870.	
13.2 Light as a particle	872
Evidence for light as a particle, 873.—How much light is one photon?, 875.—Wave-particle duality, 880.—Nonlocality and entanglement, 883.—Photons in three dimensions, 889.	
13.3 Matter as a wave	891
Electrons as waves, 892.—Dispersive waves, 896.—Bound states, 899.—The uncertainty principle, 902.—Electrons in electric fields, 904.—The Schrödinger equation, 906.	

13.4 The atom	920
Classifying states, 920.—Three dimensions, 923.—Quantum numbers, 925.—The hydrogen atom, 927.—Energies of states in hydrogen, 929.—Electron spin, 936.—Atoms with more than one electron, 939.	
Problems	942
Exercises	954
14 Additional Topics in Quantum Physics	
14.1 The Stern-Gerlach experiment	959
14.2 Rotation and vibration.	962
Types of excitations, 962.—Vibration, 962.—Rotation, 963.—Corrections to semiclassical energies, 964.	
14.3 ★A tiny bit of linear algebra.	967
14.4 The underlying structure of quantum mechanics, part 1.	969
The time-dependent Schrödinger equation, 969.—Unitarity, 972.	
14.5 Methods for solving the Schrödinger equation.	973
Cut-and-paste solutions, 973.—Separability, 976.	
14.6 The underlying structure of quantum mechanics, part 2.	977
Observables, 977.—The inner product, 984.—Completeness, 989.—The Schrödinger equation in general, 991.—Summary of the structure of quantum mechanics, 993.	
14.7 Applications to the two-state system.	993
A proton in a magnetic field, 993.—The ammonia molecule, 995.	
14.8 Energy-time uncertainty	998
Classical uncertainty relations, 998.—Energy-time uncertainty, 998.	
14.9 Randomization of phase.	1001
Randomization of phase in a measurement, 1001.—Decoherence, 1002.	
14.10 Quantum computing and the no-cloning theorem	1004
14.11 More about entanglement	1007
Problems	1011
Three essential mathematical skills	1018
Programming with python	1023
Appendix 2: Miscellany	1025
Appendix 3: Photo Credits	1031
Appendix 4: Useful Data	1070
Notation and terminology, compared with other books, 1070.—Notation and units, 1071.—Fundamental constants, 1071.—Metric prefixes, 1072.—Nonmetric units, 1072.—The Greek alphabet, 1072.—Subatomic particles, 1072.—Earth, moon, and sun, 1073.—The periodic table, 1073.—Atomic masses, 1073.	
Appendix 5: Summary	1074



The Mars Climate Orbiter is prepared for its mission. The laws of physics are the same everywhere, even on Mars, so the probe could be designed based on the laws of physics as discovered on earth. There is unfortunately another reason why this spacecraft is relevant to the topics of this chapter: it was destroyed attempting to enter Mars' atmosphere because engineers at Lockheed Martin forgot to convert data on engine thrusts from pounds into the metric unit of force (newtons) before giving the information to NASA. Conversions are important!

Chapter 0

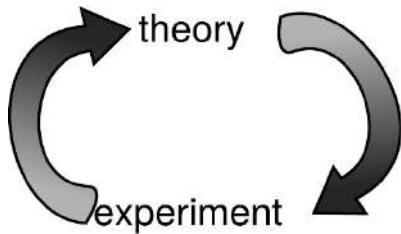
Introduction and Review

0.1 Introduction and review

If you drop your shoe and a coin side by side, they hit the ground at the same time. Why doesn't the shoe get there first, since gravity is pulling harder on it? How does the lens of your eye work, and why do your eye's muscles need to squash its lens into different shapes in order to focus on objects nearby or far away? These are the kinds of questions that physics tries to answer about the behavior of light and matter, the two things that the universe is made of.

0.1.1 The scientific method

Until very recently in history, no progress was made in answering questions like these. Worse than that, the *wrong* answers written by thinkers like the ancient Greek physicist Aristotle were accepted without question for thousands of years. Why is it that scientific knowledge has progressed more since the Renaissance than it had in all the preceding millennia since the beginning of recorded history? Undoubtedly the industrial revolution is part of the answer.



a / Science is a cycle of theory and experiment.

Building its centerpiece, the steam engine, required improved techniques for precise construction and measurement. (Early on, it was considered a major advance when English machine shops learned to build pistons and cylinders that fit together with a gap narrower than the thickness of a penny.) But even before the industrial revolution, the pace of discovery had picked up, mainly because of the introduction of the modern scientific method. Although it evolved over time, most scientists today would agree on something like the following list of the basic principles of the scientific method:

(1) *Science is a cycle of theory and experiment.* Scientific theories are created to explain the results of experiments that were created under certain conditions. A successful theory will also make new predictions about new experiments under new conditions. Eventually, though, it always seems to happen that a new experiment comes along, showing that under certain conditions the theory is not a good approximation or is not valid at all. The ball is then back in the theorists' court. If an experiment disagrees with the current theory, the theory has to be changed, not the experiment.

(2) *Theories should both predict and explain.* The requirement of predictive power means that a theory is only meaningful if it predicts something that can be checked against experimental measurements that the theorist did not already have at hand. That is, a theory should be testable. Explanatory value means that many phenomena should be accounted for with few basic principles. If you answer every “why” question with “because that’s the way it is,” then your theory has no explanatory value. Collecting lots of data without being able to find any basic underlying principles is not science.

(3) *Experiments should be reproducible.* An experiment should be treated with suspicion if it only works for one person, or only in one part of the world. Anyone with the necessary skills and equipment should be able to get the same results from the same experiment. This implies that science transcends national and ethnic boundaries; you can be sure that nobody is doing actual science who claims that their work is “Aryan, not Jewish,” “Marxist, not bourgeois,” or “Christian, not atheistic.” An experiment cannot be reproduced if it is secret, so science is necessarily a public enterprise.



b / A satirical drawing of an alchemist's laboratory. H. Cock, after a drawing by Peter Brueghel the Elder (16th century).

As an example of the cycle of theory and experiment, a vital step toward modern chemistry was the experimental observation that the chemical elements could not be transformed into each other, e.g., lead could not be turned into gold. This led to the theory that chemical reactions consisted of rearrangements of the elements in different combinations, without any change in the identities of the elements themselves. The theory worked for hundreds of years, and was confirmed experimentally over a wide range of pressures and

temperatures and with many combinations of elements. Only in the twentieth century did we learn that one element could be transformed into one another under the conditions of extremely high pressure and temperature existing in a nuclear bomb or inside a star. That observation didn't completely invalidate the original theory of the immutability of the elements, but it showed that it was only an approximation, valid at ordinary temperatures and pressures.

self-check A

A psychic conducts seances in which the spirits of the dead speak to the participants. He says he has special psychic powers not possessed by other people, which allow him to "channel" the communications with the spirits. What part of the scientific method is being violated here?

▷ Answer, p. 1057

The scientific method as described here is an idealization, and should not be understood as a set procedure for doing science. Scientists have as many weaknesses and character flaws as any other group, and it is very common for scientists to try to discredit other people's experiments when the results run contrary to their own favored point of view. Successful science also has more to do with luck, intuition, and creativity than most people realize, and the restrictions of the scientific method do not stifle individuality and self-expression any more than the fugue and sonata forms stifled Bach and Haydn. There is a recent tendency among social scientists to go even further and to deny that the scientific method even exists, claiming that science is no more than an arbitrary social system that determines what ideas to accept based on an in-group's criteria. I think that's going too far. If science is an arbitrary social ritual, it would seem difficult to explain its effectiveness in building such useful items as airplanes, CD players, and sewers. If alchemy and astrology were no less scientific in their methods than chemistry and astronomy, what was it that kept them from producing anything useful?

Discussion Questions

Consider whether or not the scientific method is being applied in the following examples. If the scientific method is not being applied, are the people whose actions are being described performing a useful human activity, albeit an unscientific one?

A Acupuncture is a traditional medical technique of Asian origin in which small needles are inserted in the patient's body to relieve pain. Many doctors trained in the west consider acupuncture unworthy of experimental study because if it had therapeutic effects, such effects could not be explained by their theories of the nervous system. Who is being more scientific, the western or eastern practitioners?

B Goethe, a German poet, is less well known for his theory of color. He published a book on the subject, in which he argued that scientific apparatus for measuring and quantifying color, such as prisms, lenses and colored filters, could not give us full insight into the ultimate meaning of color, for instance the cold feeling evoked by blue and green or the heroic sentiments inspired by red. Was his work scientific?

C A child asks why things fall down, and an adult answers “because of gravity.” The ancient Greek philosopher Aristotle explained that rocks fell because it was their nature to seek out their natural place, in contact with the earth. Are these explanations scientific?

D Buddhism is partly a psychological explanation of human suffering, and psychology is of course a science. The Buddha could be said to have engaged in a cycle of theory and experiment, since he worked by trial and error, and even late in his life he asked his followers to challenge his ideas. Buddhism could also be considered reproducible, since the Buddha told his followers they could find enlightenment for themselves if they followed a certain course of study and discipline. Is Buddhism a scientific pursuit?

0.1.2 What is physics?

Given for one instant an intelligence which could comprehend all the forces by which nature is animated and the respective positions of the things which compose it...nothing would be uncertain, and the future as the past would be laid out before its eyes.

Pierre Simon de Laplace

Physics is the use of the scientific method to find out the basic principles governing light and matter, and to discover the implications of those laws. Part of what distinguishes the modern outlook from the ancient mind-set is the assumption that there are rules by which the universe functions, and that those laws can be at least partially understood by humans. From the Age of Reason through the nineteenth century, many scientists began to be convinced that the laws of nature not only could be known but, as claimed by Laplace, those laws could in principle be used to predict everything about the universe’s future if complete information was available about the present state of all light and matter. In subsequent sections, I’ll describe two general types of limitations on prediction using the laws of physics, which were only recognized in the twentieth century.

Matter can be defined as anything that is affected by gravity, i.e., that has weight or would have weight if it was near the Earth or another star or planet massive enough to produce measurable gravity. Light can be defined as anything that can travel from one place to another through empty space and can influence matter, but has no weight. For example, sunlight can influence your body by heating it or by damaging your DNA and giving you skin cancer. The physicist’s definition of light includes a variety of phenomena that are not visible to the eye, including radio waves, microwaves, x-rays, and gamma rays. These are the “colors” of light that do not

happen to fall within the narrow violet-to-red range of the rainbow that we can see.

self-check B

At the turn of the 20th century, a strange new phenomenon was discovered in vacuum tubes: mysterious rays of unknown origin and nature. These rays are the same as the ones that shoot from the back of your TV's picture tube and hit the front to make the picture. Physicists in 1895 didn't have the faintest idea what the rays were, so they simply named them "cathode rays," after the name for the electrical contact from which they sprang. A fierce debate raged, complete with nationalistic overtones, over whether the rays were a form of light or of matter. What would they have had to do in order to settle the issue? ▷

Answer, p. 1057

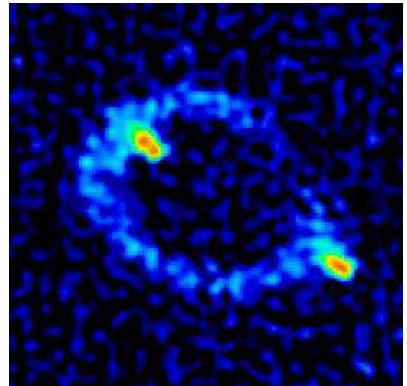
Many physical phenomena are not themselves light or matter, but are properties of light or matter or interactions between light and matter. For instance, motion is a property of all light and some matter, but it is not itself light or matter. The pressure that keeps a bicycle tire blown up is an interaction between the air and the tire. Pressure is not a form of matter in and of itself. It is as much a property of the tire as of the air. Analogously, sisterhood and employment are relationships among people but are not people themselves.

Some things that appear weightless actually do have weight, and so qualify as matter. Air has weight, and is thus a form of matter even though a cubic inch of air weighs less than a grain of sand. A helium balloon has weight, but is kept from falling by the force of the surrounding more dense air, which pushes up on it. Astronauts in orbit around the Earth have weight, and are falling along a curved arc, but they are moving so fast that the curved arc of their fall is broad enough to carry them all the way around the Earth in a circle. They perceive themselves as being weightless because their space capsule is falling along with them, and the floor therefore does not push up on their feet.

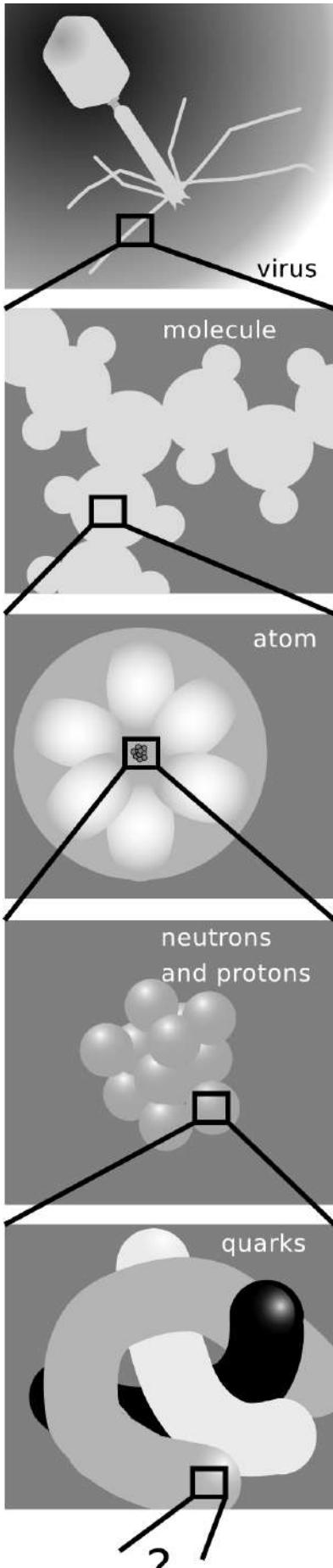
Optional Topic: Modern Changes in the Definition of Light and Matter

Einstein predicted as a consequence of his theory of relativity that light would after all be affected by gravity, although the effect would be extremely weak under normal conditions. His prediction was borne out by observations of the bending of light rays from stars as they passed close to the sun on their way to the Earth. Einstein's theory also implied the existence of black holes, stars so massive and compact that their intense gravity would not even allow light to escape. (These days there is strong evidence that black holes exist.)

Einstein's interpretation was that light doesn't really have mass, but that energy is affected by gravity just like mass is. The energy in a light beam is equivalent to a certain amount of mass, given by the famous equation $E = mc^2$, where c is the speed of light. Because the speed



c / This telescope picture shows two images of the same distant object, an exotic, very luminous object called a quasar. This is interpreted as evidence that a massive, dark object, possibly a black hole, happens to be between us and it. Light rays that would otherwise have missed the earth on either side have been bent by the dark object's gravity so that they reach us. The actual direction to the quasar is presumably in the center of the image, but the light along that central line doesn't get to us because it is absorbed by the dark object. The quasar is known by its catalog number, MG1131+0456, or more informally as Einstein's Ring.



of light is such a big number, a large amount of energy is equivalent to only a very small amount of mass, so the gravitational force on a light ray can be ignored for most practical purposes.

There is however a more satisfactory and fundamental distinction between light and matter, which should be understandable to you if you have had a chemistry course. In chemistry, one learns that electrons obey the Pauli exclusion principle, which forbids more than one electron from occupying the same orbital if they have the same spin. The Pauli exclusion principle is obeyed by the subatomic particles of which matter is composed, but disobeyed by the particles, called photons, of which a beam of light is made.

Einstein's theory of relativity is discussed more fully in book 6 of this series.

The boundary between physics and the other sciences is not always clear. For instance, chemists study atoms and molecules, which are what matter is built from, and there are some scientists who would be equally willing to call themselves physical chemists or chemical physicists. It might seem that the distinction between physics and biology would be clearer, since physics seems to deal with inanimate objects. In fact, almost all physicists would agree that the basic laws of physics that apply to molecules in a test tube work equally well for the combination of molecules that constitutes a bacterium. (Some might believe that something more happens in the minds of humans, or even those of cats and dogs.) What differentiates physics from biology is that many of the scientific theories that describe living things, while ultimately resulting from the fundamental laws of physics, cannot be rigorously derived from physical principles.

Isolated systems and reductionism

To avoid having to study everything at once, scientists isolate the things they are trying to study. For instance, a physicist who wants to study the motion of a rotating gyroscope would probably prefer that it be isolated from vibrations and air currents. Even in biology, where field work is indispensable for understanding how living things relate to their entire environment, it is interesting to note the vital historical role played by Darwin's study of the Galápagos Islands, which were conveniently isolated from the rest of the world. Any part of the universe that is considered apart from the rest can be called a "system."

Physics has had some of its greatest successes by carrying this process of isolation to extremes, subdividing the universe into smaller and smaller parts. Matter can be divided into atoms, and the behavior of individual atoms can be studied. Atoms can be split apart into their constituent neutrons, protons and electrons. Protons and neutrons appear to be made out of even smaller particles called quarks, and there have even been some claims of experimental ev-

idence that quarks have smaller parts inside them. This method of splitting things into smaller and smaller parts and studying how those parts influence each other is called reductionism. The hope is that the seemingly complex rules governing the larger units can be better understood in terms of simpler rules governing the smaller units. To appreciate what reductionism has done for science, it is only necessary to examine a 19th-century chemistry textbook. At that time, the existence of atoms was still doubted by some, electrons were not even suspected to exist, and almost nothing was understood of what basic rules governed the way atoms interacted with each other in chemical reactions. Students had to memorize long lists of chemicals and their reactions, and there was no way to understand any of it systematically. Today, the student only needs to remember a small set of rules about how atoms interact, for instance that atoms of one element cannot be converted into another via chemical reactions, or that atoms from the right side of the periodic table tend to form strong bonds with atoms from the left side.

Discussion Questions

- A** I've suggested replacing the ordinary dictionary definition of light with a more technical, more precise one that involves weightlessness. It's still possible, though, that the stuff a lightbulb makes, ordinarily called "light," does have some small amount of weight. Suggest an experiment to attempt to measure whether it does.
- B** Heat is weightless (i.e., an object becomes no heavier when heated), and can travel across an empty room from the fireplace to your skin, where it influences you by heating you. Should heat therefore be considered a form of light by our definition? Why or why not?
- C** Similarly, should sound be considered a form of light?

0.1.3 How to learn physics

For as knowledges are now delivered, there is a kind of contract of error between the deliverer and the receiver; for he that delivereth knowledge desireth to deliver it in such a form as may be best believed, and not as may be best examined; and he that receiveth knowledge desireth rather present satisfaction than expectant inquiry.

Francis Bacon

Many students approach a science course with the idea that they can succeed by memorizing the formulas, so that when a problem is assigned on the homework or an exam, they will be able to plug numbers in to the formula and get a numerical result on their calculator. Wrong! That's not what learning science is about! There is a big difference between memorizing formulas and understanding concepts. To start with, different formulas may apply in different situations. One equation might represent a definition, which is always true. Another might be a very specific equation for the speed

of an object sliding down an inclined plane, which would not be true if the object was a rock drifting down to the bottom of the ocean. If you don't work to understand physics on a conceptual level, you won't know which formulas can be used when.

Most students taking college science courses for the first time also have very little experience with interpreting the meaning of an equation. Consider the equation $w = A/h$ relating the width of a rectangle to its height and area. A student who has not developed skill at interpretation might view this as yet another equation to memorize and plug in to when needed. A slightly more savvy student might realize that it is simply the familiar formula $A = wh$ in a different form. When asked whether a rectangle would have a greater or smaller width than another with the same area but a smaller height, the unsophisticated student might be at a loss, not having any numbers to plug in on a calculator. The more experienced student would know how to reason about an equation involving division — if h is smaller, and A stays the same, then w must be bigger. Often, students fail to recognize a sequence of equations as a derivation leading to a final result, so they think all the intermediate steps are equally important formulas that they should memorize.

When learning any subject at all, it is important to become as actively involved as possible, rather than trying to read through all the information quickly without thinking about it. It is a good idea to read and think about the questions posed at the end of each section of these notes as you encounter them, so that you know you have understood what you were reading.

Many students' difficulties in physics boil down mainly to difficulties with math. Suppose you feel confident that you have enough mathematical preparation to succeed in this course, but you are having trouble with a few specific things. In some areas, the brief review given in this chapter may be sufficient, but in other areas it probably will not. Once you identify the areas of math in which you are having problems, get help in those areas. Don't limp along through the whole course with a vague feeling of dread about something like scientific notation. The problem will not go away if you ignore it. The same applies to essential mathematical skills that you are learning in this course for the first time, such as vector addition.

Sometimes students tell me they keep trying to understand a certain topic in the book, and it just doesn't make sense. The worst thing you can possibly do in that situation is to keep on staring at the same page. Every textbook explains certain things badly — even mine! — so the best thing to do in this situation is to look at a different book. Instead of college textbooks aimed at the same mathematical level as the course you're taking, you may in some cases find that high school books or books at a lower math level

give clearer explanations.

Finally, when reviewing for an exam, don't simply read back over the text and your lecture notes. Instead, try to use an active method of reviewing, for instance by discussing some of the discussion questions with another student, or doing homework problems you hadn't done the first time.

0.1.4 Velocity and acceleration

Calculus was invented by a physicist, Isaac Newton, because he needed it as a tool for calculating velocity and acceleration; in your introductory calculus course, velocity and acceleration were probably presented as some of the first applications.

If an object's position as a function of time is given by the function $x(t)$, then its velocity and acceleration are given by the first and second derivatives with respect to time,

$$v = \frac{dx}{dt}$$

and

$$a = \frac{d^2 x}{dt^2}.$$

The notation relates in a logical way to the units of the quantities. Velocity has units of m/s, and that makes sense because dx is interpreted as an infinitesimally small distance, with units of meters, and dt as an infinitesimally small time, with units of seconds. The seemingly weird and inconsistent placement of the superscripted twos in the notation for the acceleration is likewise meant to suggest the units: something on top with units of meters, and something on the bottom with units of seconds squared.

Velocity and acceleration have completely different physical interpretations. Velocity is a matter of opinion. Right now as you sit in a chair and read this book, you could say that your velocity was zero, but an observer watching the Earth rotate would say that you had a velocity of hundreds of miles an hour. Acceleration represents a *change* in velocity, and it's not a matter of opinion. Accelerations produce physical effects, and don't occur unless there's a force to cause them. For example, gravitational forces on Earth cause falling objects to have an acceleration of 9.8 m/s^2 .

Constant acceleration

example 1

- ▷ How high does a diving board have to be above the water if the diver is to have as much as 1.0 s in the air?
- ▷ The diver starts at rest, and has an acceleration of 9.8 m/s^2 . We need to find a connection between the distance she travels and time it takes. In other words, we're looking for information

about the function $x(t)$, given information about the acceleration. To go from acceleration to position, we need to integrate twice:

$$\begin{aligned}x &= \int \int a dt dt \\&= \int (at + v_0) dt \quad [v_0 \text{ is a constant of integration.}] \\&= \int at dt \quad [v_0 \text{ is zero because she's dropping from rest.}] \\&= \frac{1}{2}at^2 + x_0 \quad [x_0 \text{ is a constant of integration.}] \\&= \frac{1}{2}at^2 \quad [x_0 \text{ can be zero if we define it that way.}]\end{aligned}$$

Note some of the good problem-solving habits demonstrated here. We solve the problem symbolically, and only plug in numbers at the very end, once all the algebra and calculus are done. One should also make a habit, after finding a symbolic result, of checking whether the dependence on the variables make sense. A greater value of t in this expression would lead to a greater value for x ; that makes sense, because if you want more time in the air, you're going to have to jump from higher up. A greater acceleration also leads to a greater height; this also makes sense, because the stronger gravity is, the more height you'll need in order to stay in the air for a given amount of time. Now we plug in numbers.

$$\begin{aligned}x &= \frac{1}{2} (9.8 \text{ m/s}^2) (1.0 \text{ s})^2 \\&= 4.9 \text{ m}\end{aligned}$$

Note that when we put in the numbers, we check that the units work out correctly, $(\text{m/s}^2)(\text{s})^2 = \text{m}$. We should also check that the result makes sense: 4.9 meters is pretty high, but not unreasonable.

The notation dq in calculus represents an infinitesimally small change in the variable q . The corresponding notation for a finite change in a variable is Δq . For example, if q represents the value of a certain stock on the stock market, and the value falls from $q_0 = 5$ dollars initially to $q_f = 3$ dollars finally, then $\Delta q = -2$ dollars. When we study linear functions, whose slopes are constant, the derivative is synonymous with the slope of the line, and dy/dx is the same thing as $\Delta y/\Delta x$, the rise over the run.

Under conditions of constant acceleration, we can relate velocity and time,

$$a = \frac{\Delta v}{\Delta t},$$

or, as in the example 1, position and time,

$$x = \frac{1}{2}at^2 + v_0t + x_0.$$

It can also be handy to have a relation involving velocity and position, eliminating time. Straightforward algebra gives

$$v_f^2 = v_o^2 + 2a\Delta x,$$

where v_f is the final velocity, v_o the initial velocity, and Δx the distance traveled.

▷ *Solved problem: Dropping a rock on Mars* page 49, problem 17

▷ *Solved problem: The Dodge Viper* page 50, problem 19

0.1.5 Self-evaluation

The introductory part of a book like this is hard to write, because every student arrives at this starting point with a different preparation. One student may have grown up outside the U.S. and so may be completely comfortable with the metric system, but may have had an algebra course in which the instructor passed too quickly over scientific notation. Another student may have already taken vector calculus, but may have never learned the metric system. The following self-evaluation is a checklist to help you figure out what you need to study to be prepared for the rest of the course.

If you disagree with this statement...	you should study this section:
I am familiar with the basic metric units of meters, kilograms, and seconds, and the most common metric prefixes: milli- (m), kilo- (k), and centi- (c).	subsection 0.1.6 Basic of the Metric System
I am familiar with these less common metric prefixes: mega- (M), micro- (μ), and nano- (n).	subsection 0.1.7 Less Common Metric Prefixes
I am comfortable with scientific notation.	subsection 0.1.8 Scientific Notation
I can confidently do metric conversions.	subsection 0.1.9 Conversions
I understand the purpose and use of significant figures.	subsection 0.1.10 Significant Figures

It wouldn't hurt you to skim the sections you think you already know about, and to do the self-checks in those sections.

0.1.6 Basics of the metric system

The metric system

Every country in the world besides the U.S. uses a system of units known in English as the “metric system.¹” This system is entirely decimal, thanks to the same eminently logical people who brought about the French Revolution. In deference to France, the system’s official name is the Système International, or SI, meaning International System. The system uses a single, consistent set of Greek and Latin prefixes that modify the basic units. Each prefix stands for a power of ten, and has an abbreviation that can be combined with the symbol for the unit. For instance, the meter is a unit of distance. The prefix kilo- stands for 10^3 , so a kilometer, 1 km, is a thousand meters.

The basic units of the SI are the meter for distance, the second for time, and the kilogram (not the gram) for mass.

The following are the most common metric prefixes. You should memorize them.

prefix	meaning	example
kilo-	10^3	60 kg = a person’s mass
centi-	10^{-2}	28 cm = height of a piece of paper
milli-	10^{-3}	1 ms = time for one vibration of a guitar string playing the note D

The prefix centi-, meaning 10^{-2} , is only used in the centimeter; a hundredth of a gram would not be written as 1 cg but as 10 mg. The centi- prefix can be easily remembered because a cent is 10^{-2} dollars. The official SI abbreviation for seconds is “s” (not “sec”) and grams are “g” (not “gm”).

The second

When I stated briefly above that the second was a unit of time, it may not have occurred to you that this was not much of a definition. We can make a dictionary-style definition of a term like “time,” or give a general description like Isaac Newton’s: “Absolute, true, and mathematical time, of itself, and from its own nature, flows equably without relation to anything external...” Newton’s characterization sounds impressive, but physicists today would consider it useless as a definition of time. Today, the physical sciences are based on operational definitions, which means definitions that spell out the actual steps (operations) required to measure something numerically.

In an era when our toasters, pens, and coffee pots tell us the time, it is far from obvious to most people what is the fundamental operational definition of time. Until recently, the hour, minute, and second were defined operationally in terms of the time required for

¹Liberia and Myanmar have not legally adopted metric units, but use them in everyday life.

the earth to rotate about its axis. Unfortunately, the Earth's rotation is slowing down slightly, and by 1967 this was becoming an issue in scientific experiments requiring precise time measurements. The second was therefore redefined as the time required for a certain number of vibrations of the light waves emitted by a cesium atoms in a lamp constructed like a familiar neon sign but with the neon replaced by cesium. The new definition not only promises to stay constant indefinitely, but for scientists is a more convenient way of calibrating a clock than having to carry out astronomical measurements.

self-check C

What is a possible operational definition of how strong a person is? ▶

Answer, p. 1057

The meter

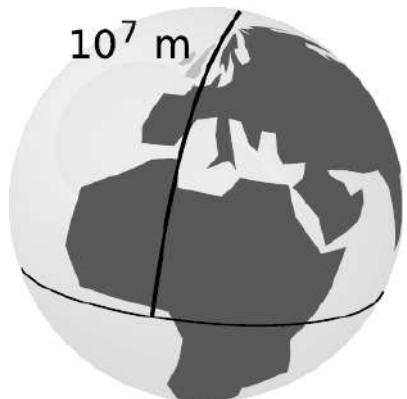
The French originally defined the meter as 10^{-7} times the distance from the equator to the north pole, as measured through Paris (of course). Even if the definition was operational, the operation of traveling to the north pole and laying a surveying chain behind you was not one that most working scientists wanted to carry out. Fairly soon, a standard was created in the form of a metal bar with two scratches on it. This was replaced by an atomic standard in 1960, and finally in 1983 by the current definition, which is that the speed of light has a defined value in units of m/s.

The kilogram

The third base unit of the SI is the kilogram, a unit of mass. Mass is intended to be a measure of the amount of a substance, but that is not an operational definition. Bathroom scales work by measuring our planet's gravitational attraction for the object being weighed, but using that type of scale to define mass operationally would be undesirable because gravity varies in strength from place to place on the earth. The kilogram was for a long time defined by a physical artifact (figure f), but in 2019 it was redefined by giving a defined value to Planck's constant (p. 877), which plays a fundamental role in the description of the atomic world.

Combinations of metric units

Just about anything you want to measure can be measured with some combination of meters, kilograms, and seconds. Speed can be measured in m/s, volume in m^3 , and density in kg/m^3 . Part of what makes the SI great is this basic simplicity. No more funny units like a cord of wood, a bolt of cloth, or a jigger of whiskey. No more liquid and dry measure. Just a simple, consistent set of units. The SI measures put together from meters, kilograms, and seconds make up the mks system. For example, the mks unit of speed is m/s, not km/hr.



e / The original definition of the meter.



f / A duplicate of the Paris kilogram, maintained at the Danish National Metrology Institute. As of 2019, the kilogram is no longer defined in terms of a physical standard.

Checking units

A useful technique for finding mistakes in one's algebra is to analyze the units associated with the variables.

Checking units

example 2

▷ Jae starts from the formula $V = \frac{1}{3}Ah$ for the volume of a cone, where A is the area of its base, and h is its height. He wants to find an equation that will tell him how tall a conical tent has to be in order to have a certain volume, given its radius. His algebra goes like this:

$$\begin{aligned}[1] V &= \frac{1}{3}Ah \\ [2] A &= \pi r^2 \\ [3] V &= \frac{1}{3}\pi r^2 h \\ [4] h &= \frac{\pi r^2}{3V} \end{aligned}$$

Is his algebra correct? If not, find the mistake.

▷ Line 4 is supposed to be an equation for the height, so the units of the expression on the right-hand side had better equal meters. The pi and the 3 are unitless, so we can ignore them. In terms of units, line 4 becomes

$$m = \frac{m^2}{m^3} = \frac{1}{m}.$$

This is false, so there must be a mistake in the algebra. The units of lines 1, 2, and 3 check out, so the mistake must be in the step from line 3 to line 4. In fact the result should have been

$$h = \frac{3V}{\pi r^2}.$$

Now the units check: $m = m^3/m^2$.

Discussion Question

A Isaac Newton wrote, "...the natural days are truly unequal, though they are commonly considered as equal, and used for a measure of time... It may be that there is no such thing as an equable motion, whereby time may be accurately measured. All motions may be accelerated or retarded..." Newton was right. Even the modern definition of the second in terms of light emitted by cesium atoms is subject to variation. For instance, magnetic fields could cause the cesium atoms to emit light with a slightly different rate of vibration. What makes us think, though, that a pendulum clock is more accurate than a sundial, or that a cesium atom is a more accurate timekeeper than a pendulum clock? That is, how can one test experimentally how the accuracies of different time standards compare?

0.1.7 Less common metric prefixes

The following are three metric prefixes which, while less common than the ones discussed previously, are well worth memorizing.

prefix	meaning	example
mega-	M 10^6	6.4 Mm = radius of the earth
micro-	μ 10^{-6}	10 μm = size of a white blood cell
nano-	n 10^{-9}	0.154 nm = distance between carbon nuclei in an ethane molecule

Note that the abbreviation for micro is the Greek letter mu, μ — a common mistake is to confuse it with m (milli) or M (mega).

There are other prefixes even less common, used for extremely large and small quantities. For instance, 1 femtometer = 10^{-15} m is a convenient unit of distance in nuclear physics, and 1 gigabyte = 10^9 bytes is used for computers' hard disks. The international committee that makes decisions about the SI has recently even added some new prefixes that sound like jokes, e.g., 1 yoctogram = 10^{-24} g is about half the mass of a proton. In the immediate future, however, you're unlikely to see prefixes like "yocto-" and "zepto-" used except perhaps in trivia contests at science-fiction conventions or other geekfests.

self-check D

Suppose you could slow down time so that according to your perception, a beam of light would move across a room at the speed of a slow walk. If you perceived a nanosecond as if it was a second, how would you perceive a microsecond?

▷ Answer, p. 1057

0.1.8 Scientific notation

Most of the interesting phenomena in our universe are not on the human scale. It would take about 1,000,000,000,000,000,000 bacteria to equal the mass of a human body. When the physicist Thomas Young discovered that light was a wave, it was back in the bad old days before scientific notation, and he was obliged to write that the time required for one vibration of the wave was 1/500 of a millionth of a millionth of a second. Scientific notation is a less awkward way to write very large and very small numbers such as these. Here's a quick review.

Scientific notation means writing a number in terms of a product of something from 1 to 10 and something else that is a power of ten. For instance,

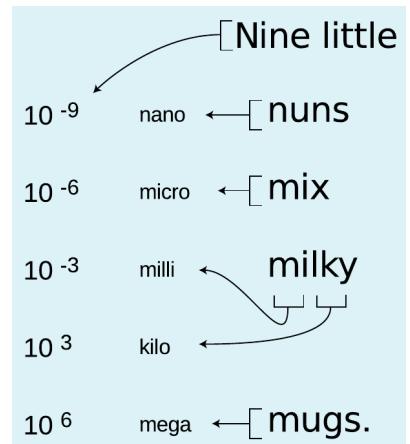
$$32 = 3.2 \times 10^1$$

$$320 = 3.2 \times 10^2$$

$$3200 = 3.2 \times 10^3 \dots$$

Each number is ten times bigger than the previous one.

Since 10^1 is ten times smaller than 10^2 , it makes sense to use



g / This is a mnemonic to help you remember the most important metric prefixes. The word "little" is to remind you that the list starts with the prefixes used for small quantities and builds upward. The exponent changes by 3, except that of course that we do not need a special prefix for 10^0 , which equals one.

the notation 10^0 to stand for one, the number that is in turn ten times smaller than 10^1 . Continuing on, we can write 10^{-1} to stand for 0.1, the number ten times smaller than 10^0 . Negative exponents are used for small numbers:

$$\begin{aligned}3.2 &= 3.2 \times 10^0 \\0.32 &= 3.2 \times 10^{-1} \\0.032 &= 3.2 \times 10^{-2} \quad \dots\end{aligned}$$

A common source of confusion is the notation used on the displays of many calculators. Examples:

3.2×10^6	(written notation)
3.2E+6	(notation on some calculators)
3.2^6	(notation on some other calculators)

The last example is particularly unfortunate, because 3.2^6 really stands for the number $3.2 \times 3.2 \times 3.2 \times 3.2 \times 3.2 \times 3.2 = 1074$, a totally different number from $3.2 \times 10^6 = 3200000$. The calculator notation should never be used in writing. It's just a way for the manufacturer to save money by making a simpler display.

self-check E

A student learns that 10^4 bacteria, standing in line to register for classes at Paramecium Community College, would form a queue of this size:



The student concludes that 10^2 bacteria would form a line of this length:



Why is the student incorrect?

▷ Answer, p. 1057

0.1.9 Conversions

Conversions are one of the three essential mathematical skills, summarized on pp.1018-1020, that you need for success in this course.

I suggest you avoid memorizing lots of conversion factors between SI units and U.S. units, but two that do come in handy are:

$$1 \text{ inch} = 2.54 \text{ cm}$$

An object with a weight on Earth of 2.2 pounds-force has a mass of 1 kg.

The first one is the present definition of the inch, so it's exact. The second one is not exact, but is good enough for most purposes. (U.S. units of force and mass are confusing, so it's a good thing they're not used in science. In U.S. units, the unit of force is the pound-

force, and the best unit to use for mass is the slug, which is about 14.6 kg.)

More important than memorizing conversion factors is understanding the right method for doing conversions. Even within the SI, you may need to convert, say, from grams to kilograms. Different people have different ways of thinking about conversions, but the method I'll describe here is systematic and easy to understand. The idea is that if 1 kg and 1000 g represent the same mass, then we can consider a fraction like

$$\frac{10^3 \text{ g}}{1 \text{ kg}}$$

to be a way of expressing the number one. This may bother you. For instance, if you type 1000/1 into your calculator, you will get 1000, not one. Again, different people have different ways of thinking about it, but the justification is that it helps us to do conversions, and it works! Now if we want to convert 0.7 kg to units of grams, we can multiply kg by the number one:

$$0.7 \text{ kg} \times \frac{10^3 \text{ g}}{1 \text{ kg}}$$

If you're willing to treat symbols such as "kg" as if they were variables as used in algebra (which they're really not), you can then cancel the kg on top with the kg on the bottom, resulting in

$$0.7 \cancel{\text{kg}} \times \frac{10^3 \text{ g}}{1 \cancel{\text{kg}}} = 700 \text{ g.}$$

To convert grams to kilograms, you would simply flip the fraction upside down.

One advantage of this method is that it can easily be applied to a series of conversions. For instance, to convert one year to units of seconds,

$$1 \text{ year} \times \frac{365 \cancel{\text{days}}}{1 \cancel{\text{year}}} \times \frac{24 \cancel{\text{hours}}}{1 \cancel{\text{day}}} \times \frac{60 \cancel{\text{min}}}{1 \cancel{\text{hour}}} \times \frac{60 \cancel{\text{s}}}{1 \cancel{\text{min}}} = 3.15 \times 10^7 \text{ s.}$$

Should that exponent be positive, or negative?

A common mistake is to write the conversion fraction incorrectly. For instance the fraction

$$\frac{10^3 \text{ kg}}{1 \text{ g}} \quad (\text{incorrect})$$

does not equal one, because 10^3 kg is the mass of a car, and 1 g is the mass of a raisin. One correct way of setting up the conversion

factor would be

$$\frac{10^{-3} \text{ kg}}{1 \text{ g}} \quad (\text{correct}).$$

You can usually detect such a mistake if you take the time to check your answer and see if it is reasonable.

If common sense doesn't rule out either a positive or a negative exponent, here's another way to make sure you get it right. There are big prefixes and small prefixes:

big prefixes: k M
small prefixes: m μ n

(It's not hard to keep straight which are which, since "mega" and "micro" are evocative, and it's easy to remember that a kilometer is bigger than a meter and a millimeter is smaller.) In the example above, we want the top of the fraction to be the same as the bottom. Since k is a big prefix, we need to *compensate* by putting a small number like 10^{-3} in front of it, not a big number like 10^3 .

- ▷ *Solved problem: a simple conversion* page 47, problem 6
- ▷ *Solved problem: the geometric mean* page 47, problem 8

Discussion Question

A Each of the following conversions contains an error. In each case, explain what the error is.

- (a) $1000 \text{ kg} \times \frac{1 \text{ kg}}{1000 \text{ g}} = 1 \text{ g}$
- (b) $50 \text{ m} \times \frac{1 \text{ cm}}{100 \text{ m}} = 0.5 \text{ cm}$
- (c) "Nano" is 10^{-9} , so there are 10^{-9} nm in a meter.
- (d) "Micro" is 10^{-6} , so 1 kg is $10^6 \mu\text{g}$.

0.1.10 Significant figures

The international governing body for football ("soccer" in the US) says the ball should have a circumference of 68 to 70 cm. Taking the middle of this range and dividing by π gives a diameter of approximately 21.96338214668155633610595934540698196 cm. The digits after the first few are completely meaningless. Since the circumference could have varied by about a centimeter in either direction, the diameter is fuzzy by something like a third of a centimeter. We say that the additional, random digits are not *significant figures*. If you write down a number with a lot of gratuitous insignificant figures, it shows a lack of scientific literacy and imples to other people a greater precision than you really have.

As a rule of thumb, the result of a calculation has as many significant figures, or "sig figs," as the least accurate piece of data that went in. In the example with the soccer ball, it didn't do us any good to know π to dozens of digits, because the bottleneck in the precision of the result was the figure for the circumference, which

was two sig figs. The result is 22 cm. The rule of thumb works best for multiplication and division.

For calculations involving multiplication and division, a given fractional or “percent” error in one of the inputs causes the same fractional error in the output. The number of digits in a number provides a rough measure of its possible fractional error. These are called significant figures or “sig figs.” Examples:

3.14	3 sig figs
3.1	2 sig figs
0.03	1 sig fig, because the zeroes are just placeholders
3.0×10^1	2 sig figs
30	could be 1 or 2 sig figs, since we can't tell if the 0 is a placeholder or a real sig fig

In such calculations, your result should not have more than the number of sig figs in the least accurate piece of data you started with.

- Sig figs in the area of a triangle* *example 3*
- ▷ A triangle has an area of 6.45 m^2 and a base with a width of 4.0138 m. Find its height.
 - ▷ The area is related to the base and height by $A = bh/2$.

$$\begin{aligned} h &= \frac{2A}{b} \\ &= 3.21391200358762 \text{ m} \quad (\text{calculator output}) \\ &= 3.21 \text{ m} \end{aligned}$$

The given data were 3 sig figs and 5 sig figs. We’re limited by the less accurate piece of data, so the final result is 3 sig figs. The additional digits on the calculator don’t mean anything, and if we communicated them to another person, we would create the false impression of having determined h with more precision than we really obtained.

self-check F

The following quote is taken from an editorial by Norimitsu Onishi in the New York Times, August 18, 2002.

Consider Nigeria. Everyone agrees it is Africa’s most populous nation. But what is its population? The United Nations says 114 million; the State Department, 120 million. The World Bank says 126.9 million, while the Central Intelligence Agency puts it at 126,635,626.

What should bother you about this?

▷ Answer, p. 1058

Dealing correctly with significant figures can save you time! Of-

ten, students copy down numbers from their calculators with eight significant figures of precision, then type them back in for a later calculation. That's a waste of time, unless your original data had that kind of incredible precision.

self-check G

How many significant figures are there in each of the following measurements?

- (1) 9.937 m
- (2) 4.0 s
- (3) 0.0000000000000037 kg

▷ Answer, p. 1058

The rules about significant figures are only rules of thumb, and are not a substitute for careful thinking. For instance, \$20.00 + \$0.05 is \$20.05. It need not and should not be rounded off to \$20. In general, the sig fig rules work best for multiplication and division, and we sometimes also apply them when doing a complicated calculation that involves many types of operations. For simple addition and subtraction, it makes more sense to maintain a fixed number of digits after the decimal point.

When in doubt, don't use the sig fig rules at all. Instead, intentionally change one piece of your initial data by the maximum amount by which you think it could have been off, and recalculate the final result. The digits on the end that are completely reshuffled are the ones that are meaningless, and should be omitted.

A nonlinear function

example 4

- ▷ How many sig figs are there in $\sin 88.7^\circ$?

▷ We're using a sine function, which isn't addition, subtraction, multiplication, or division. It would be reasonable to guess that since the input angle had 3 sig figs, so would the output. But if this was an important calculation and we really needed to know, we would do the following:

$$\sin 88.7^\circ = 0.999742609322698$$

$$\sin 88.8^\circ = 0.999780683474846$$

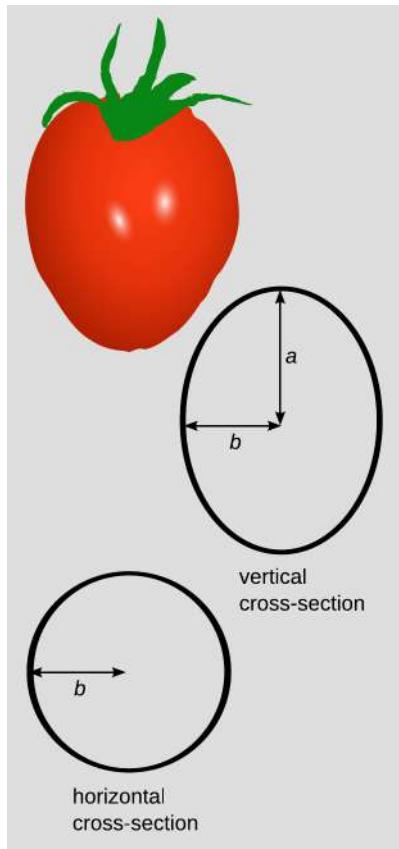
Surprisingly, the result appears to have as many as 5 sig figs, not just 3:

$$\sin 88.7^\circ = 0.99974,$$

where the final 4 is uncertain but may have some significance. The unexpectedly high precision of the result is because the sine function is nearing its maximum at 90 degrees, where the graph flattens out and becomes insensitive to the input angle.

0.1.11 A note about diagrams

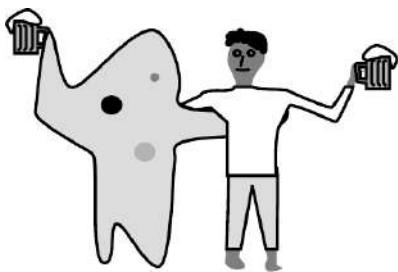
A quick note about diagrams. Often when you solve a problem, the best way to get started and organize your thoughts is by drawing a diagram. For an artist, it's desirable to be able to draw a



h / A diagram of a tomato.

recognizable, realistic, perspective picture of a tomato, like the one at the top of figure h. But in science and engineering, we usually don't draw solid figures in perspective, because that would make it difficult to label distances and angles. Usually we want views or cross-sections that project the object into its planes of symmetry, as in the line drawings in the figure.

0.2 Scaling and order-of-magnitude estimates



a / Amoebas this size are seldom encountered.

0.2.1 Introduction

Why can't an insect be the size of a dog? Some skinny stretched-out cells in your spinal cord are a meter tall — why does nature display no single cells that are not just a meter tall, but a meter wide, and a meter thick as well? Believe it or not, these are questions that can be answered fairly easily without knowing much more about physics than you already do. The only mathematical technique you really need is the humble conversion, applied to area and volume.

Area and volume

Area can be defined by saying that we can copy the shape of interest onto graph paper with $1\text{ cm} \times 1\text{ cm}$ squares and count the number of squares inside. Fractions of squares can be estimated by eye. We then say the area equals the number of squares, in units of square cm. Although this might seem less "pure" than computing areas using formulae like $A = \pi r^2$ for a circle or $A = wh/2$ for a triangle, those formulae are not useful as definitions of area because they cannot be applied to irregularly shaped areas.

Units of square cm are more commonly written as cm^2 in science. Of course, the unit of measurement symbolized by "cm" is not an algebra symbol standing for a number that can be literally multiplied by itself. But it is advantageous to write the units of area that way and treat the units as if they were algebra symbols. For instance, if you have a rectangle with an area of 6 m^2 and a width of 2 m, then calculating its length as $(6\text{ m}^2)/(2\text{ m}) = 3\text{ m}$ gives a result that makes sense both numerically and in terms of units. This algebra-style treatment of the units also ensures that our methods of converting units work out correctly. For instance, if we accept the fraction

$$\frac{100\text{ cm}}{1\text{ m}}$$

as a valid way of writing the number one, then one times one equals one, so we should also say that one can be represented by

$$\frac{100\text{ cm}}{1\text{ m}} \times \frac{100\text{ cm}}{1\text{ m}},$$

which is the same as

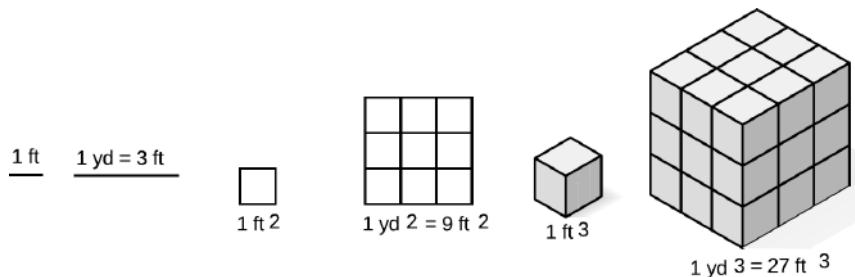
$$\frac{10000\text{ cm}^2}{1\text{ m}^2}.$$

That means the conversion factor from square meters to square centimeters is a factor of 10^4 , i.e., a square meter has 10^4 square centimeters in it.

All of the above can be easily applied to volume as well, using one-cubic-centimeter blocks instead of squares on graph paper.

To many people, it seems hard to believe that a square meter equals 10000 square centimeters, or that a cubic meter equals a

million cubic centimeters — they think it would make more sense if there were 100 cm^2 in 1 m^2 , and 100 cm^3 in 1 m^3 , but that would be incorrect. The examples shown in figure b aim to make the correct answer more believable, using the traditional U.S. units of feet and yards. (One foot is 12 inches, and one yard is three feet.)



b / Visualizing conversions of area and volume using traditional U.S. units.

self-check H

Based on figure b, convince yourself that there are 9 ft^2 in a square yard, and 27 ft^3 in a cubic yard, then demonstrate the same thing symbolically (i.e., with the method using fractions that equal one). \triangleright Answer, p. 1058

\triangleright Solved problem: converting mm^2 to cm^2 page 51, problem 31

\triangleright Solved problem: scaling a liter page 52, problem 40

Discussion Question

A How many square centimeters are there in a square inch? (1 inch = 2.54 cm) First find an approximate answer by making a drawing, then derive the conversion factor more accurately using the symbolic method.

0.2.2 Scaling of area and volume

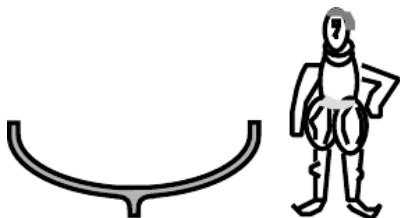
Great fleas have lesser fleas
Upon their backs to bite 'em.
And lesser fleas have lesser still,
And so ad infinitum.

Jonathan Swift

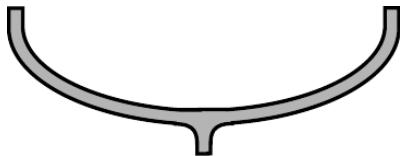
Now how do these conversions of area and volume relate to the questions I posed about sizes of living things? Well, imagine that you are shrunk like Alice in Wonderland to the size of an insect. One way of thinking about the change of scale is that what used to look like a centimeter now looks like perhaps a meter to you, because you're so much smaller. If area and volume scaled according to most people's intuitive, incorrect expectations, with 1 m^2 being the same as 100 cm^2 , then there would be no particular reason why nature should behave any differently on your new, reduced scale. But nature does behave differently now that you're small. For instance, you will find that you can walk on water, and jump



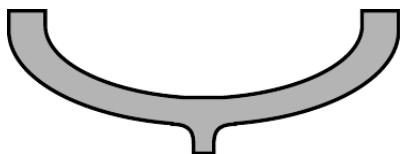
c / Galileo Galilei (1564-1642).



d / The small boat holds up just fine.



e / A larger boat built with the same proportions as the small one will collapse under its own weight.



f / A boat this large needs to have timbers that are thicker compared to its size.

to many times your own height. The physicist Galileo Galilei had the basic insight that the scaling of area and volume determines how natural phenomena behave differently on different scales. He first reasoned about mechanical structures, but later extended his insights to living things, taking the then-radical point of view that at the fundamental level, a living organism should follow the same laws of nature as a machine. We will follow his lead by first discussing machines and then living things.

Galileo on the behavior of nature on large and small scales

One of the world's most famous pieces of scientific writing is Galileo's Dialogues Concerning the Two New Sciences. Galileo was an entertaining writer who wanted to explain things clearly to laypeople, and he livened up his work by casting it in the form of a dialogue among three people. Salviati is really Galileo's alter ego. Simplicio is the stupid character, and one of the reasons Galileo got in trouble with the Church was that there were rumors that Simplicio represented the Pope. Sagredo is the earnest and intelligent student, with whom the reader is supposed to identify. (The following excerpts are from the 1914 translation by Crew and de Salvio.)

SAGREDO: Yes, that is what I mean; and I refer especially to his last assertion which I have always regarded as false... ; namely, that in speaking of these and other similar machines one cannot argue from the small to the large, because many devices which succeed on a small scale do not work on a large scale. Now, since mechanics has its foundations in geometry, where mere size [is unimportant], I do not see that the properties of circles, triangles, cylinders, cones and other solid figures will change with their size. If, therefore, a large machine be constructed in such a way that its parts bear to one another the same ratio as in a smaller one, and if the smaller is sufficiently strong for the purpose for which it is designed, I do not see why the larger should not be able to withstand any severe and destructive tests to which it may be subjected.

Salviati contradicts Sagredo:

SALVIATI: ...Please observe, gentlemen, how facts which at first seem improbable will, even on scant explanation, drop the cloak which has hidden them and stand forth in naked and simple beauty. Who does not know that a horse falling from a height of three or four cubits will break his bones, while a dog falling from the same height or a cat from a height of eight or ten cubits will suffer no injury? Equally harmless would be the fall of a grasshopper from a tower or the fall of an ant from the distance of the moon.

The point Galileo is making here is that small things are sturdier

in proportion to their size. There are a lot of objections that could be raised, however. After all, what does it really mean for something to be “strong”, to be “strong in proportion to its size,” or to be strong “out of proportion to its size?” Galileo hasn’t given operational definitions of things like “strength,” i.e., definitions that spell out how to measure them numerically.

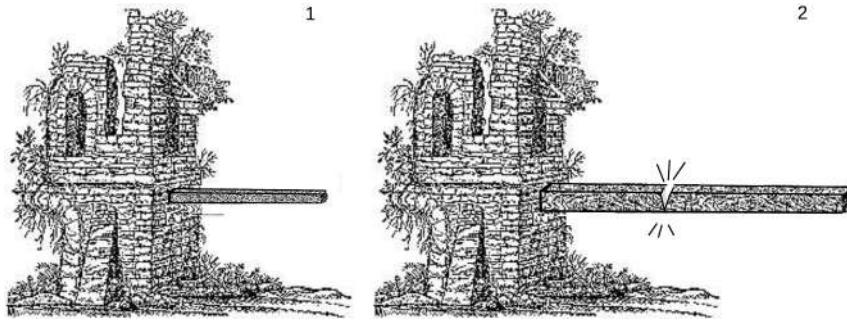
Also, a cat is shaped differently from a horse — an enlarged photograph of a cat would not be mistaken for a horse, even if the photo-doctoring experts at the National Inquirer made it look like a person was riding on its back. A grasshopper is not even a mammal, and it has an exoskeleton instead of an internal skeleton. The whole argument would be a lot more convincing if we could do some isolation of variables, a scientific term that means to change only one thing at a time, isolating it from the other variables that might have an effect. If size is the variable whose effect we’re interested in seeing, then we don’t really want to compare things that are different in size but also different in other ways.

SALVIATI: ...we asked the reason why [shipbuilders] employed stocks, scaffolding, and bracing of larger dimensions for launching a big vessel than they do for a small one; and [an old man] answered that they did this in order to avoid the danger of the ship parting under its own heavy weight, a danger to which small boats are not subject?

After this entertaining but not scientifically rigorous beginning, Galileo starts to do something worthwhile by modern standards. He simplifies everything by considering the strength of a wooden plank. The variables involved can then be narrowed down to the type of wood, the width, the thickness, and the length. He also gives an operational definition of what it means for the plank to have a certain strength “in proportion to its size,” by introducing the concept of a plank that is the longest one that would not snap under its own weight if supported at one end. If you increased its length by the slightest amount, without increasing its width or thickness, it would break. He says that if one plank is the same shape as another but a different size, appearing like a reduced or enlarged photograph of the other, then the planks would be strong “in proportion to their sizes” if both were just barely able to support their own weight.



h / Galileo discusses planks made of wood, but the concept may be easier to imagine with clay. All three clay rods in the figure were originally the same shape. The medium-sized one was twice the height, twice the length, and twice the width of the small one, and similarly the large one was twice as big as the medium one in all its linear dimensions. The big one has four times the linear dimensions of the small one, 16 times the cross-sectional area when cut perpendicular to the page, and 64 times the volume. That means that the big one has 64 times the weight to support, but only 16 times the strength compared to the smallest one.



g / 1. This plank is as long as it can be without collapsing under its own weight. If it was a hundredth of an inch longer, it would collapse.
 2. This plank is made out of the same kind of wood. It is twice as thick, twice as long, and twice as wide. It will collapse under its own weight.

Also, Galileo is doing something that would be frowned on in modern science: he is mixing experiments whose results he has actually observed (building boats of different sizes), with experiments that he could not possibly have done (dropping an ant from the height of the moon). He now relates how he has done actual experiments with such planks, and found that, according to this operational definition, they are not strong in proportion to their sizes. The larger one breaks. He makes sure to tell the reader how important the result is, via Sagredo's astonished response:

SAGREDO: My brain already reels. My mind, like a cloud momentarily illuminated by a lightning flash, is for an instant filled with an unusual light, which now beckons to me and which now suddenly mingles and obscures strange, crude ideas. From what you have said it appears to me impossible to build two similar structures of the same material, but of different sizes and have them proportionately strong.

In other words, this specific experiment, using things like wooden planks that have no intrinsic scientific interest, has very wide implications because it points out a general principle, that nature acts differently on different scales.

To finish the discussion, Galileo gives an explanation. He says that the strength of a plank (defined as, say, the weight of the heaviest boulder you could put on the end without breaking it) is proportional to its cross-sectional area, that is, the surface area of the fresh wood that would be exposed if you sawed through it in the middle. Its weight, however, is proportional to its volume.²

How do the volume and cross-sectional area of the longer plank compare with those of the shorter plank? We have already seen,

²Galileo makes a slightly more complicated argument, taking into account the effect of leverage (torque). The result I'm referring to comes out the same regardless of this effect.

while discussing conversions of the units of area and volume, that these quantities don't act the way most people naively expect. You might think that the volume and area of the longer plank would both be doubled compared to the shorter plank, so they would increase in proportion to each other, and the longer plank would be equally able to support its weight. You would be wrong, but Galileo knows that this is a common misconception, so he has Salviati address the point specifically:

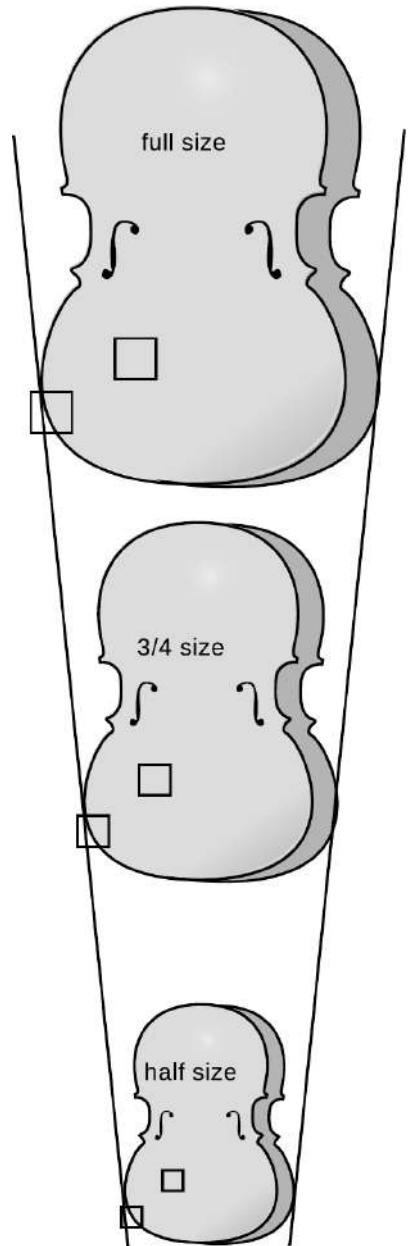
SALVIATI: ... Take, for example, a cube two inches on a side so that each face has an area of four square inches and the total area, i.e., the sum of the six faces, amounts to twenty-four square inches; now imagine this cube to be sawed through three times [with cuts in three perpendicular planes] so as to divide it into eight smaller cubes, each one inch on the side, each face one inch square, and the total surface of each cube six square inches instead of twenty-four in the case of the larger cube. It is evident therefore, that the surface of the little cube is only one-fourth that of the larger, namely, the ratio of six to twenty-four; but the volume of the solid cube itself is only one-eighth; the volume, and hence also the weight, diminishes therefore much more rapidly than the surface... You see, therefore, Simplicio, that I was not mistaken when ... I said that the surface of a small solid is comparatively greater than that of a large one.

The same reasoning applies to the planks. Even though they are not cubes, the large one could be sawed into eight small ones, each with half the length, half the thickness, and half the width. The small plank, therefore, has more surface area in proportion to its weight, and is therefore able to support its own weight while the large one breaks.

Scaling of area and volume for irregularly shaped objects

You probably are not going to believe Galileo's claim that this has deep implications for all of nature unless you can be convinced that the same is true for any shape. Every drawing you've seen so far has been of squares, rectangles, and rectangular solids. Clearly the reasoning about sawing things up into smaller pieces would not prove anything about, say, an egg, which cannot be cut up into eight smaller egg-shaped objects with half the length.

Is it always true that something half the size has one quarter the surface area and one eighth the volume, even if it has an irregular shape? Take the example of a child's violin. Violins are made for small children in smaller size to accomodate their small bodies. Figure i shows a full-size violin, along with two violins made with half and 3/4 of the normal length.³ Let's study the surface area of



i / The area of a shape is proportional to the square of its linear dimensions, even if the shape is irregular.

³The customary terms "half-size" and "3/4-size" actually don't describe the

the front panels of the three violins.

Consider the square in the interior of the panel of the full-size violin. In the 3/4-size violin, its height and width are both smaller by a factor of 3/4, so the area of the corresponding, smaller square becomes $3/4 \times 3/4 = 9/16$ of the original area, not 3/4 of the original area. Similarly, the corresponding square on the smallest violin has half the height and half the width of the original one, so its area is 1/4 the original area, not half.

The same reasoning works for parts of the panel near the edge, such as the part that only partially fills in the other square. The entire square scales down the same as a square in the interior, and in each violin the same fraction (about 70%) of the square is full, so the contribution of this part to the total area scales down just the same.

Since any small square region or any small region covering part of a square scales down like a square object, the entire surface area of an irregularly shaped object changes in the same manner as the surface area of a square: scaling it down by 3/4 reduces the area by a factor of 9/16, and so on.

In general, we can see that any time there are two objects with the same shape, but different linear dimensions (i.e., one looks like a reduced photo of the other), the ratio of their areas equals the ratio of the squares of their linear dimensions:

$$\frac{A_1}{A_2} = \left(\frac{L_1}{L_2} \right)^2.$$

Note that it doesn't matter where we choose to measure the linear size, L , of an object. In the case of the violins, for instance, it could have been measured vertically, horizontally, diagonally, or even from the bottom of the left f-hole to the middle of the right f-hole. We just have to measure it in a consistent way on each violin. Since all the parts are assumed to shrink or expand in the same manner, the ratio L_1/L_2 is independent of the choice of measurement.

It is also important to realize that it is completely unnecessary to have a formula for the area of a violin. It is only possible to derive simple formulas for the areas of certain shapes like circles, rectangles, triangles and so on, but that is no impediment to the type of reasoning we are using.

Sometimes it is inconvenient to write all the equations in terms of ratios, especially when more than two objects are being compared. A more compact way of rewriting the previous equation is

$$A \propto L^2.$$

sizes in any accurate way. They're really just standard, arbitrary marketing labels.

The symbol “ \propto ” means “is proportional to.” Scientists and engineers often speak about such relationships verbally using the phrases “scales like” or “goes like,” for instance “area goes like length squared.”

All of the above reasoning works just as well in the case of volume. Volume goes like length cubed:

$$V \propto L^3.$$

self-check 1

When a car or truck travels over a road, there is wear and tear on the road surface, which incurs a cost. Studies show that the cost C per kilometer of travel is related to the weight per axle w by $C \propto w^4$. Translate this into a statement about ratios. \triangleright Answer, p. 1058

If different objects are made of the same material with the same density, $\rho = m/V$, then their masses, $m = \rho V$, are proportional to L^3 . (The symbol for density is ρ , the lower-case Greek letter “rho.”)

An important point is that all of the above reasoning about scaling only applies to objects that are the same shape. For instance, a piece of paper is larger than a pencil, but has a much greater surface-to-volume ratio.

Scaling of the area of a triangle

example 5

- \triangleright In figure k, the larger triangle has sides twice as long. How many times greater is its area?

Correct solution #1: Area scales in proportion to the square of the linear dimensions, so the larger triangle has four times more area ($2^2 = 4$).

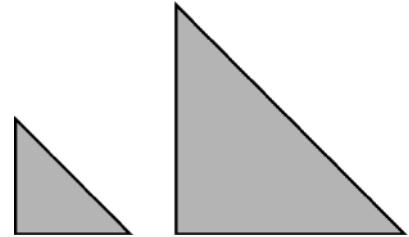
Correct solution #2: You could cut the larger triangle into four of the smaller size, as shown in fig. (b), so its area is four times greater. (This solution is correct, but it would not work for a shape like a circle, which can't be cut up into smaller circles.)

Correct solution #3: The area of a triangle is given by

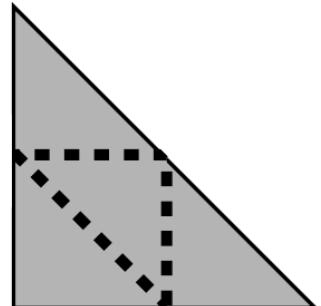
$A = bh/2$, where b is the base and h is the height. The areas of the triangles are

$$\begin{aligned} A_1 &= b_1 h_1 / 2 \\ A_2 &= b_2 h_2 / 2 \\ &= (2b_1)(2h_1) / 2 \\ &= 2b_1 h_1 \\ A_2/A_1 &= (2b_1 h_1) / (b_1 h_1 / 2) \\ &= 4 \end{aligned}$$

(Although this solution is correct, it is a lot more work than solution #1, and it can only be used in this case because a triangle is a simple geometric shape, and we happen to know a formula for its area.)



k / Example 5. The big triangle has four times more area than the little one.



b / A tricky way of solving example 5, explained in solution #2.

Correct solution #4: The area of a triangle is $A = bh/2$. The comparison of the areas will come out the same as long as the ratios of the linear sizes of the triangles is as specified, so let's just say $b_1 = 1.00 \text{ m}$ and $b_2 = 2.00 \text{ m}$. The heights are then also $h_1 = 1.00 \text{ m}$ and $h_2 = 2.00 \text{ m}$, giving areas $A_1 = 0.50 \text{ m}^2$ and $A_2 = 2.00 \text{ m}^2$, so $A_2/A_1 = 4.00$.

(The solution is correct, but it wouldn't work with a shape for whose area we don't have a formula. Also, the numerical calculation might make the answer of 4.00 appear inexact, whereas solution #1 makes it clear that it is exactly 4.)

Incorrect solution: The area of a triangle is $A = bh/2$, and if you plug in $b = 2.00 \text{ m}$ and $h = 2.00 \text{ m}$, you get $A = 2.00 \text{ m}^2$, so the bigger triangle has 2.00 times more area. (This solution is incorrect because no comparison has been made with the smaller triangle.)

Scaling of the volume of a sphere

example 6

- ▷ In figure m, the larger sphere has a radius that is five times greater. How many times greater is its volume?

Correct solution #1: Volume scales like the third power of the linear size, so the larger sphere has a volume that is 125 times greater ($5^3 = 125$).

Correct solution #2: The volume of a sphere is $V = (4/3)\pi r^3$, so

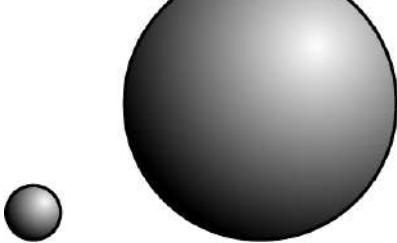
$$V_1 = \frac{4}{3}\pi r_1^3$$

$$V_2 = \frac{4}{3}\pi r_2^3$$

$$= \frac{4}{3}\pi(5r_1)^3$$

$$= \frac{500}{3}\pi r_1^3$$

$$V_2/V_1 = \left(\frac{500}{3}\pi r_1^3\right) / \left(\frac{4}{3}\pi r_1^3\right) = 125$$



m / Example 6. The big sphere has 125 times more volume than the little one.



n / Example 7. The 48-point "S" has 1.78 times more area than the 36-point "s."

$$V_1 = \frac{4}{3}\pi r_1^3$$

$$V_2 = \frac{4}{3}\pi r_2^3$$

$$= \frac{4}{3}\pi \cdot 5r_1^3$$

$$= \frac{20}{3}\pi r_1^3$$

$$V_2/V_1 = \left(\frac{20}{3}\pi r_1^3\right) / \left(\frac{4}{3}\pi r_1^3\right) = 5$$

(The solution is incorrect because $(5r_1)^3$ is not the same as $5r_1^3$.)

Scaling of a more complex shape

example 7

- ▷ The first letter “S” in figure n is in a 36-point font, the second in 48-point. How many times more ink is required to make the larger “S”? (Points are a unit of length used in typography.)

Correct solution: The amount of ink depends on the area to be covered with ink, and area is proportional to the square of the linear dimensions, so the amount of ink required for the second “S” is greater by a factor of $(48/36)^2 = 1.78$.

Incorrect solution: The length of the curve of the second “S” is longer by a factor of $48/36 = 1.33$, so 1.33 times more ink is required.

(The solution is wrong because it assumes incorrectly that the width of the curve is the same in both cases. Actually both the width and the length of the curve are greater by a factor of $48/36$, so the area is greater by a factor of $(48/36)^2 = 1.78$.)

Reasoning about ratios and proportionalities is one of the three essential mathematical skills, summarized on pp.1018-1020, that you need for success in this course.

▷ *Solved problem: a telescope gathers light* page 51, problem 32

▷ *Solved problem: distance from an earthquake* page 51, problem 33

Discussion Questions

A A toy fire engine is 1/30 the size of the real one, but is constructed from the same metal with the same proportions. How many times smaller is its weight? How many times less red paint would be needed to paint it?

B Galileo spends a lot of time in his dialog discussing what really happens when things break. He discusses everything in terms of Aristotle’s now-discredited explanation that things are hard to break, because if something breaks, there has to be a gap between the two halves with nothing in between, at least initially. Nature, according to Aristotle, “abhors a vacuum,” i.e., nature doesn’t “like” empty space to exist. Of course, air will rush into the gap immediately, but at the very moment of breaking, Aristotle imagined a vacuum in the gap. Is Aristotle’s explanation of why it is hard to break things an experimentally testable statement? If so, how could it be tested experimentally?

0.2.3 Order-of-magnitude estimates

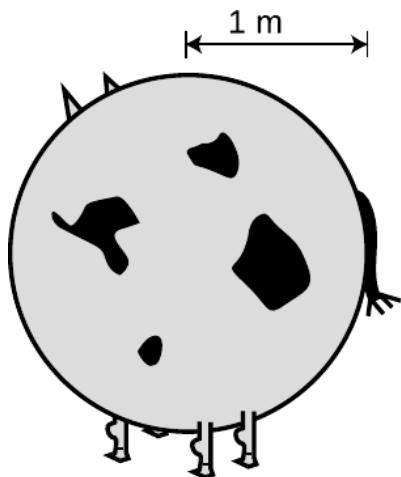
It is the mark of an instructed mind to rest satisfied with the degree of precision that the nature of the subject permits and not to seek an exactness where only an approximation of the truth is possible.

Aristotle

It is a common misconception that science must be exact. For instance, in the Star Trek TV series, it would often happen that Captain Kirk would ask Mr. Spock, “Spock, we’re in a pretty bad situation. What do you think are our chances of getting out of



o / Can you guess how many jelly beans are in the jar? If you try to guess directly, you will almost certainly underestimate. The right way to do it is to estimate the linear dimensions, then get the volume indirectly. See problem 44, p. 53.



p / Consider a spherical cow.

here?" The scientific Mr. Spock would answer with something like, "Captain, I estimate the odds as 237.345 to one." In reality, he could not have estimated the odds with six significant figures of accuracy, but nevertheless one of the hallmarks of a person with a good education in science is the ability to make estimates that are likely to be at least somewhere in the right ballpark. In many such situations, it is often only necessary to get an answer that is off by no more than a factor of ten in either direction. Since things that differ by a factor of ten are said to differ by one order of magnitude, such an estimate is called an order-of-magnitude estimate. The tilde, \sim , is used to indicate that things are only of the same order of magnitude, but not exactly equal, as in

$$\text{odds of survival} \sim 100 \text{ to one.}$$

The tilde can also be used in front of an individual number to emphasize that the number is only of the right order of magnitude.

Although making order-of-magnitude estimates seems simple and natural to experienced scientists, it's a mode of reasoning that is completely unfamiliar to most college students. Some of the typical mental steps can be illustrated in the following example.

Cost of transporting tomatoes (incorrect solution) example 8

- ▷ Roughly what percentage of the price of a tomato comes from the cost of transporting it in a truck?
- ▷ The following incorrect solution illustrates one of the main ways you can go wrong in order-of-magnitude estimates.

Incorrect solution: Let's say the trucker needs to make a \$400 profit on the trip. Taking into account her benefits, the cost of gas, and maintenance and payments on the truck, let's say the total cost is more like \$2000. I'd guess about 5000 tomatoes would fit in the back of the truck, so the extra cost per tomato is 40 cents. That means the cost of transporting one tomato is comparable to the cost of the tomato itself. Transportation really adds a lot to the cost of produce, I guess.

The problem is that the human brain is not very good at estimating area or volume, so it turns out the estimate of 5000 tomatoes fitting in the truck is way off. That's why people have a hard time at those contests where you are supposed to estimate the number of jellybeans in a big jar. Another example is that most people think their families use about 10 gallons of water per day, but in reality the average is about 300 gallons per day. When estimating area or volume, you are much better off estimating linear dimensions, and computing volume from the linear dimensions. Here's a better solution to the problem about the tomato truck:

Cost of transporting tomatoes (correct solution) example 9
As in the previous solution, say the cost of the trip is \$2000. The

dimensions of the bin are probably $4\text{ m} \times 2\text{ m} \times 1\text{ m}$, for a volume of 8 m^3 . Since the whole thing is just an order-of-magnitude estimate, let's round that off to the nearest power of ten, 10 m^3 . The shape of a tomato is complicated, and I don't know any formula for the volume of a tomato shape, but since this is just an estimate, let's pretend that a tomato is a cube, $0.05\text{ m} \times 0.05\text{ m} \times 0.05\text{ m}$, for a volume of $1.25 \times 10^{-4}\text{ m}^3$. Since this is just a rough estimate, let's round that to 10^{-4} m^3 . We can find the total number of tomatoes by dividing the volume of the bin by the volume of one tomato: $10\text{ m}^3 / 10^{-4}\text{ m}^3 = 10^5$ tomatoes. The transportation cost per tomato is $\$2000 / 10^5$ tomatoes = $\$0.02/\text{tomato}$. That means that transportation really doesn't contribute very much to the cost of a tomato.

Approximating the shape of a tomato as a cube is an example of another general strategy for making order-of-magnitude estimates. A similar situation would occur if you were trying to estimate how many m^2 of leather could be produced from a herd of ten thousand cattle. There is no point in trying to take into account the shape of the cows' bodies. A reasonable plan of attack might be to consider a spherical cow. Probably a cow has roughly the same surface area as a sphere with a radius of about 1 m, which would be $4\pi(1\text{ m})^2$. Using the well-known facts that pi equals three, and four times three equals about ten, we can guess that a cow has a surface area of about 10 m^2 , so the herd as a whole might yield 10^5 m^2 of leather.

Estimating mass indirectly

example 10

Usually the best way to estimate mass is to estimate linear dimensions, then use those to infer volume, and then get the mass based on the volume. For example, *Amphicoelias*, shown in the figure, may have been the largest land animal ever to live. Fossils tell us the linear dimensions of an animal, but we can only indirectly guess its mass. Given the length scale in the figure, let's estimate the mass of an *Amphicoelias*.

Its torso looks like it can be approximated by a rectangular box with dimensions $10\text{ m} \times 5\text{ m} \times 3\text{ m}$, giving about $2 \times 10^2\text{ m}^3$. Living things are mostly made of water, so we assume the animal to have the density of water, 1 g/cm^3 , which converts to 10^3 kg/m^3 . This gives a mass of about $2 \times 10^5\text{ kg}$, or 200 metric tons.



The following list summarizes the strategies for getting a good order-of-magnitude estimate.

1. Don't even attempt more than one significant figure of precision.
2. Don't guess area, volume, or mass directly. Guess linear dimensions and get area, volume, or mass from them.
3. When dealing with areas or volumes of objects with complex shapes, idealize them as if they were some simpler shape, a cube or a sphere, for example.
4. Check your final answer to see if it is reasonable. If you estimate that a herd of ten thousand cattle would yield 0.01 m^2 of leather, then you have probably made a mistake with conversion factors somewhere.

Problems

The symbols ✓, ■, etc. are explained on page 53.

- 1** Correct use of a calculator: (a) Calculate $\frac{74658}{53222+97554}$ on a calculator. [Self-check: The most common mistake results in 97555.40.] ✓

- (b) Which would be more like the price of a TV, and which would be more like the price of a house, $\$3.5 \times 10^5$ or $\$3.5^5$? ■

- 2** Compute the following things. If they don't make sense because of units, say so.

- (a) 3 cm + 5 cm
(b) 1.11 m + 22 cm
(c) 120 miles + 2.0 hours
(d) 120 miles / 2.0 hours



- 3** Your backyard has brick walls on both ends. You measure a distance of 23.4 m from the inside of one wall to the inside of the other. Each wall is 29.4 cm thick. How far is it from the outside of one wall to the outside of the other? Pay attention to significant figures. ■

- 4** The speed of light is 3.0×10^8 m/s. Convert this to furlongs per fortnight. A furlong is 220 yards, and a fortnight is 14 days. An inch is 2.54 cm. ✓



- 5** Express each of the following quantities in micrograms:
(a) 10 mg, (b) 10^4 g, (c) 10 kg, (d) 100×10^3 g, (e) 1000 ng. ✓



- 6** Convert 134 mg to units of kg, writing your answer in scientific notation. ▷ Solution, p. 1037 ■

- 7** In the last century, the average age of the onset of puberty for girls has decreased by several years. Urban folklore has it that this is because of hormones fed to beef cattle, but it is more likely to be because modern girls have more body fat on the average and possibly because of estrogen-mimicking chemicals in the environment from the breakdown of pesticides. A hamburger from a hormone-implanted steer has about 0.2 ng of estrogen (about double the amount of natural beef). A serving of peas contains about 300 ng of estrogen. An adult woman produces about 0.5 mg of estrogen per day (note the different unit!). (a) How many hamburgers would a girl have to eat in one day to consume as much estrogen as an adult woman's daily production? (b) How many servings of peas?



- 8** The usual definition of the mean (average) of two numbers a and b is $(a+b)/2$. This is called the arithmetic mean. The geometric

mean, however, is defined as $(ab)^{1/2}$ (i.e., the square root of ab). For the sake of definiteness, let's say both numbers have units of mass. (a) Compute the arithmetic mean of two numbers that have units of grams. Then convert the numbers to units of kilograms and recompute their mean. Is the answer consistent? (b) Do the same for the geometric mean. (c) If a and b both have units of grams, what should we call the units of ab ? Does your answer make sense when you take the square root? (d) Suppose someone proposes to you a third kind of mean, called the superduper mean, defined as $(ab)^{1/3}$. Is this reasonable? \triangleright Solution, p. 1037 ■

9 In an article on the SARS epidemic, the May 7, 2003 New York Times discusses conflicting estimates of the disease's incubation period (the average time that elapses from infection to the first symptoms). "The study estimated it to be 6.4 days. But other statistical calculations ... showed that the incubation period could be as long as 14.22 days." What's wrong here? ■

Problem 10.

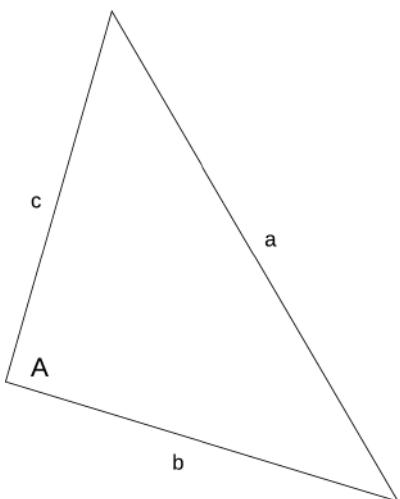
10 The photo shows the corner of a bag of pretzels. What's wrong here? ■

11 The distance to the horizon is given by the expression $\sqrt{2rh}$, where r is the radius of the Earth, and h is the observer's height above the Earth's surface. (This can be proved using the Pythagorean theorem.) Show that the units of this expression make sense. Don't try to prove the result, just check its units. (See example 2 on p. 26 for an example of how to do this.) ■

12 (a) Based on the definitions of the sine, cosine, and tangent, what units must they have? (b) A cute formula from trigonometry lets you find any angle of a triangle if you know the lengths of its sides. Using the notation shown in the figure, and letting $s = (a + b + c)/2$ be half the perimeter, we have

$$\tan A/2 = \sqrt{\frac{(s - b)(s - c)}{s(s - a)}}.$$

Show that the units of this equation make sense. In other words, check that the units of the right-hand side are the same as your answer to part a of the question. \triangleright Solution, p. 1037 ■



Problem 12.

13 A physics homework question asks, "If you start from rest and accelerate at 1.54 m/s^2 for 3.29 s , how far do you travel by the end of that time?" A student answers as follows:

$$1.54 \times 3.29 = 5.07 \text{ m}$$

His Aunt Wanda is good with numbers, but has never taken physics. She doesn't know the formula for the distance traveled under constant acceleration over a given amount of time, but she tells her nephew his answer cannot be right. How does she know? ■

14 You are looking into a deep well. It is dark, and you cannot see the bottom. You want to find out how deep it is, so you drop a rock in, and you hear a splash 3.0 seconds later. How deep is the well? ✓ ■

15 You take a trip in your spaceship to another star. Setting off, you increase your speed at a constant acceleration. Once you get half-way there, you start decelerating, at the same rate, so that by the time you get there, you have slowed down to zero speed. You see the tourist attractions, and then head home by the same method.

(a) Find a formula for the time, T , required for the round trip, in terms of d , the distance from our sun to the star, and a , the magnitude of the acceleration. Note that the acceleration is not constant over the whole trip, but the trip can be broken up into constant-acceleration parts.

(b) The nearest star to the Earth (other than our own sun) is Proxima Centauri, at a distance of $d = 4 \times 10^{16}$ m. Suppose you use an acceleration of $a = 10$ m/s², just enough to compensate for the lack of true gravity and make you feel comfortable. How long does the round trip take, in years?

(c) Using the same numbers for d and a , find your maximum speed. Compare this to the speed of light, which is 3.0×10^8 m/s. (Later in this course, you will learn that there are some new things going on in physics when one gets close to the speed of light, and that it is impossible to exceed the speed of light. For now, though, just use the simpler ideas you've learned so far.) ✓ ■

16 You climb half-way up a tree, and drop a rock. Then you climb to the top, and drop another rock. How many times greater is the velocity of the second rock on impact? Explain. (The answer is not two times greater.) ■

17 If the acceleration of gravity on Mars is $1/3$ that on Earth, how many times longer does it take for a rock to drop the same distance on Mars? Ignore air resistance. ➤ Solution, p. 1038 ■

18 A person is parachute jumping. During the time between when she leaps out of the plane and when she opens her chute, her altitude is given by an equation of the form

$$y = b - c \left(t + ke^{-t/k} \right),$$

where e is the base of natural logarithms, and b , c , and k are constants. Because of air resistance, her velocity does not increase at a steady rate as it would for an object falling in vacuum.

(a) What units would b , c , and k have to have for the equation to make sense?

(b) Find the person's velocity, v , as a function of time. [You will need to use the chain rule, and the fact that $d(e^x)/dx = e^x$.] ✓

(c) Use your answer from part (b) to get an interpretation of the constant c . [Hint: e^{-x} approaches zero for large values of x .]

- (d) Find the person's acceleration, a , as a function of time. ✓
(e) Use your answer from part (d) to show that if she waits long enough to open her chute, her acceleration will become very small. ■

19 In July 1999, Popular Mechanics carried out tests to find which car sold by a major auto maker could cover a quarter mile (402 meters) in the shortest time, starting from rest. Because the distance is so short, this type of test is designed mainly to favor the car with the greatest acceleration, not the greatest maximum speed (which is irrelevant to the average person). The winner was the Dodge Viper, with a time of 12.08 s. The car's top (and presumably final) speed was 118.51 miles per hour (52.98 m/s). (a) If a car, starting from rest and moving with *constant* acceleration, covers a quarter mile in this time interval, what is its acceleration? (b) What would be the final speed of a car that covered a quarter mile with the constant acceleration you found in part a? (c) Based on the discrepancy between your answer in part b and the actual final speed of the Viper, what do you conclude about how its acceleration changed over time? ▷ Solution, p. 1038 ■

20 The speed required for a low-earth orbit is 7.9×10^3 m/s. When a rocket is launched into orbit, it goes up a little at first to get above almost all of the atmosphere, but then tips over horizontally to build up to orbital speed. Suppose the horizontal acceleration is limited to $3g$ to keep from damaging the cargo (or hurting the crew, for a crewed flight). (a) What is the minimum distance the rocket must travel downrange before it reaches orbital speed? How much does it matter whether you take into account the initial eastward velocity due to the rotation of the earth? (b) Rather than a rocket ship, it might be advantageous to use a railgun design, in which the craft would be accelerated to orbital speeds along a railroad track. This has the advantage that it isn't necessary to lift a large mass of fuel, since the energy source is external. Based on your answer to part a, comment on the feasibility of this design for crewed launches from the earth's surface. ■

21 Consider the following passage from Alice in Wonderland, in which Alice has been falling for a long time down a rabbit hole:

Down, down, down. Would the fall *never* come to an end? "I wonder how many miles I've fallen by this time?" she said aloud. "I must be getting somewhere near the center of the earth. Let me see: that would be four thousand miles down, I think" (for, you see, Alice had learned several things of this sort in her lessons in the schoolroom, and though this was not a *very* good opportunity for showing off her knowledge, as there was no one to listen to her, still it was good practice to say it over)...

Alice doesn't know much physics, but let's try to calculate the amount of time it would take to fall four thousand miles, starting

from rest with an acceleration of 10 m/s^2 . This is really only a lower limit; if there really was a hole that deep, the fall would actually take a longer time than the one you calculate, both because there is air friction and because gravity gets weaker as you get deeper (at the center of the earth, g is zero, because the earth is pulling you equally in every direction at once). ✓ ■

22 How many cubic inches are there in a cubic foot? The answer is not 12. ✓ ■

23 Assume a dog's brain is twice as great in diameter as a cat's, but each animal's brain cells are the same size and their brains are the same shape. In addition to being a far better companion and much nicer to come home to, how many times more brain cells does a dog have than a cat? The answer is not 2. ■

24 The population density of Los Angeles is about 4000 people/ km^2 . That of San Francisco is about 6000 people/ km^2 . How many times farther away is the average person's nearest neighbor in LA than in San Francisco? The answer is not 1.5. ✓ ■

25 A hunting dog's nose has about 10 square inches of active surface. How is this possible, since the dog's nose is only about 1 in \times 1 in \times 1 in = 1 in 3 ? After all, 10 is greater than 1, so how can it fit? ■

26 Estimate the number of blades of grass on a football field. ■

27 In a computer memory chip, each bit of information (a 0 or a 1) is stored in a single tiny circuit etched onto the surface of a silicon chip. The circuits cover the surface of the chip like lots in a housing development. A typical chip stores 64 Mb (megabytes) of data, where a byte is 8 bits. Estimate (a) the area of each circuit, and (b) its linear size. ■

28 Suppose someone built a gigantic apartment building, measuring 10 km \times 10 km at the base. Estimate how tall the building would have to be to have space in it for the entire world's population to live. ■

29 A hamburger chain advertises that it has sold 10 billion Bongo Burgers. Estimate the total mass of feed required to raise the cows used to make the burgers. ■

30 Estimate the volume of a human body, in cm 3 . ■

31 How many cm 2 is 1 mm 2 ? ▷ Solution, p. 1038 ■

32 Compare the light-gathering powers of a 3-cm-diameter telescope and a 30-cm telescope. ▷ Solution, p. 1038 ■

33 One step on the Richter scale corresponds to a factor of 100 in terms of the energy absorbed by something on the surface of the Earth, e.g., a house. For instance, a 9.3-magnitude quake would

release 100 times more energy than an 8.3. The energy spreads out from the epicenter as a wave, and for the sake of this problem we'll assume we're dealing with seismic waves that spread out in three dimensions, so that we can visualize them as hemispheres spreading out under the surface of the earth. If a certain 7.6-magnitude earthquake and a certain 5.6-magnitude earthquake produce the same amount of vibration where I live, compare the distances from my house to the two epicenters.

▷ Solution, p. 1038 ■

34 In Europe, a piece of paper of the standard size, called A4, is a little narrower and taller than its American counterpart. The ratio of the height to the width is the square root of 2, and this has some useful properties. For instance, if you cut an A4 sheet from left to right, you get two smaller sheets that have the same proportions. You can even buy sheets of this smaller size, and they're called A5. There is a whole series of sizes related in this way, all with the same proportions. (a) Compare an A5 sheet to an A4 in terms of area and linear size. (b) The series of paper sizes starts from an A0 sheet, which has an area of one square meter. Suppose we had a series of boxes defined in a similar way: the B0 box has a volume of one cubic meter, two B1 boxes fit exactly inside an B0 box, and so on. What would be the dimensions of a B0 box? ✓ ■

35 Estimate the mass of one of the hairs in Albert Einstein's moustache, in units of kg. ■

36 According to folklore, every time you take a breath, you are inhaling some of the atoms exhaled in Caesar's last words. Is this true? If so, how many? ■

37 The Earth's surface is about 70% water. Mars's diameter is about half the Earth's, but it has no surface water. Compare the land areas of the two planets. ✓ ■

38 The traditional Martini glass is shaped like a cone with the point at the bottom. Suppose you make a Martini by pouring vermouth into the glass to a depth of 3 cm, and then adding gin to bring the depth to 6 cm. What are the proportions of gin and vermouth?

▷ Solution, p. 1038 ■

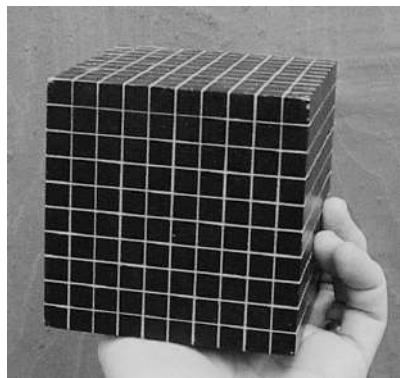
39 The central portion of a CD is taken up by the hole and some surrounding clear plastic, and this area is unavailable for storing data. The radius of the central circle is about 35% of the outer radius of the data-storing area. What percentage of the CD's area is therefore lost? ✓ ■

40 The one-liter cube in the photo has been marked off into smaller cubes, with linear dimensions one tenth those of the big one. What is the volume of each of the small cubes?

▷ Solution, p. 1038 ■



Albert Einstein, and his moustache, problem 35.



Problem 40.

41 Estimate the number of man-hours required for building the Great Wall of China. ▷ Solution, p. 1038

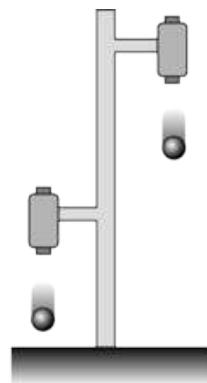
42 (a) Using the microscope photo in the figure, estimate the mass of a one cell of the *E. coli* bacterium, which is one of the most common ones in the human intestine. Note the scale at the lower right corner, which is $1 \mu\text{m}$. Each of the tubular objects in the column is one cell. (b) The feces in the human intestine are mostly bacteria (some dead, some alive), of which *E. coli* is a large and typical component. Estimate the number of bacteria in your intestines, and compare with the number of human cells in your body, which is believed to be roughly on the order of 10^{13} . (c) Interpreting your result from part b, what does this tell you about the size of a typical human cell compared to the size of a typical bacterial cell? ■



Problem 42.

43 The figure shows a practical, simple experiment for determining g to high precision. Two steel balls are suspended from electromagnets, and are released simultaneously when the electric current is shut off. They fall through unequal heights Δx_1 and Δx_2 . A computer records the sounds through a microphone as first one ball and then the other strikes the floor. From this recording, we can accurately determine the quantity T defined as $T = \Delta t_2 - \Delta t_1$, i.e., the time lag between the first and second impacts. Note that since the balls do not make any sound when they are released, we have no way of measuring the individual times Δt_2 and Δt_1 .

- (a) Find an equation for g in terms of the measured quantities T , Δx_1 and Δx_2 . ✓
- (b) Check the units of your equation.
- (c) Check that your equation gives the correct result in the case where Δx_1 is very close to zero. However, is this case realistic?
- (d) What happens when $\Delta x_1 = \Delta x_2$? Discuss this both mathematically and physically. ■



Problem 43.

44 Estimate the number of jellybeans in figure o on p. 44.

▷ Solution, p. 1039

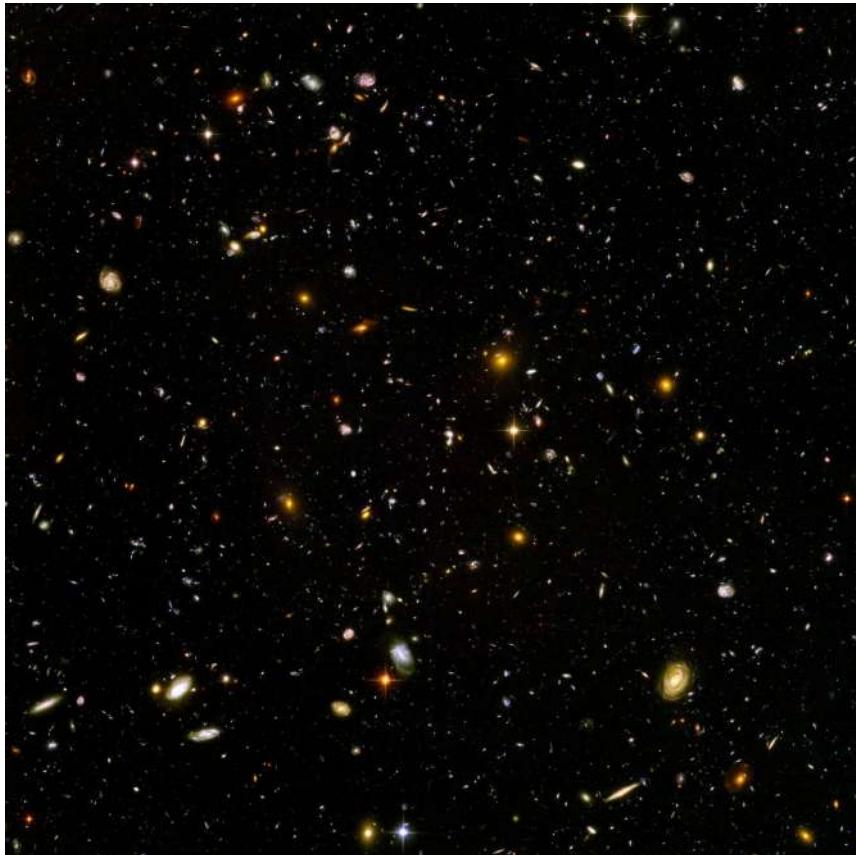
45 Let the function x be defined by $x(t) = Ae^{bt}$, where t has units of seconds and x has units of meters. (For $b < 0$, this could be a fairly accurate model of the motion of a bullet shot into a tank of oil.) Show that the Taylor series of this function makes sense if and only if A and b have certain units. ■

46 A 2002 paper by Steegmann *et al.* uses data from modern human groups like the Inuit to argue that Neanderthals in Ice Age Europe had to eat up “to 4,480 kcal per day to support strenuous winter foraging and cold resistance costs.” What’s wrong here? ■

Key to symbols:

■ easy ■ typical ■ challenging ■ difficult ■ very difficult

✓ An answer check is available at www.lightandmatter.com.



The universe has been recycling its contents ever since the Big Bang, 13.7 billion years ago.

Chapter 1

Conservation of Mass

It took just a moment for that head to fall, but a hundred years might not produce another like it.

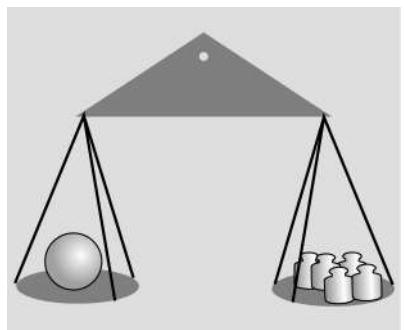
Joseph-Louis Lagrange, referring to the execution of Lavoisier on May 8, 1794

1.1 Mass

Change is impossible, claimed the ancient Greek philosopher Parmenides. His work was nonscientific, since he didn't state his ideas in a form that would allow them to be tested experimentally, but modern science nevertheless has a strong Parmenidean flavor. His main argument that change is an illusion was that something can't be turned into nothing, and likewise if you have nothing, you can't turn it into something. To make this into a scientific theory, we have to decide on a way to measure what "something" is, and we can then



a / Portrait of Monsieur Lavoisier and His Wife, by Jacques-Louis David, 1788. Lavoisier invented the concept of conservation of mass. The husband is depicted with his scientific apparatus, while in the background on the left is the portfolio belonging to Madame Lavoisier, who is thought to have been a student of David's.



b / A measurement of gravitational mass: the sphere has a gravitational mass of five kilograms.

check by measurements whether the total amount of “something” in the universe really stays constant. How much “something” is there in a rock? Does a sunbeam count as “something”? Does heat count? Motion? Thoughts and feelings?

If you look at the table of contents of this book, you’ll see that the first four chapters have the word “conservation” in them. In physics, a conservation law is a statement that the total amount of a certain physical quantity always stays the same. This chapter is about conservation of mass. The metric system is designed around a unit of distance, the meter, a unit of mass, the kilogram, and a time unit, the second. Numerical measurement of distance and time probably date back almost as far into prehistory as counting money, but mass is a more modern concept. Until scientists figured out that mass was conserved, it wasn’t obvious that there could be a single, consistent way of measuring an amount of matter, hence jiggers of whiskey and cords of wood. You may wonder why conservation of mass wasn’t discovered until relatively modern times, but it wasn’t obvious, for example, that gases had mass, and that the apparent loss of mass when wood was burned was exactly matched by the mass of the escaping gases.

Once scientists were on the track of the conservation of mass concept, they began looking for a way to define mass in terms of a definite measuring procedure. If they tried such a procedure, and the result was that it led to nonconservation of mass, then they would throw it out and try a different procedure. For instance, we might be tempted to define mass using kitchen measuring cups, i.e., as a measure of volume. Mass would then be perfectly conserved for a process like mixing marbles with peanut butter, but there would be processes like freezing water that led to a net increase in mass, and others like soaking up water with a sponge that caused a decrease. If, with the benefit of hindsight, it seems like the measuring cup definition was just plain silly, then here’s a more subtle example of a wrong definition of mass. Suppose we define it using a bathroom scale, or a more precise device such as a postal scale that works on the same principle of using gravity to compress or twist a spring. The trouble is that gravity is not equally strong all over the surface of the earth, so for instance there would be nonconservation of mass when you brought an object up to the top of a mountain, where gravity is a little weaker.

Although some of the obvious possibilities have problems, there do turn out to be at least two approaches to defining mass that lead to its being a conserved quantity, so we consider these definitions to be “right” in the pragmatic sense that what’s correct is what’s useful.

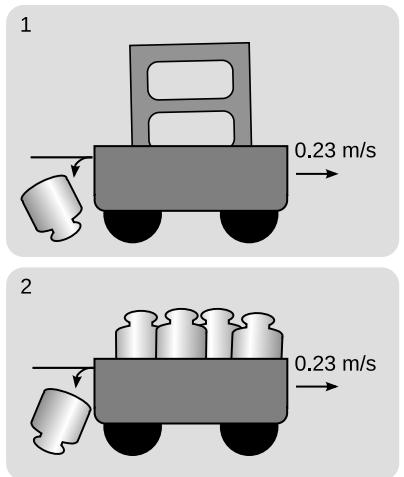
One definition that works is to use balances, but compensate for the local strength of gravity. This is the method that is used

by scientists who actually specialize in ultraprecise measurements. A standard kilogram, in the form of a platinum-iridium cylinder, is kept in a special shrine in Paris. Copies are made that balance against the standard kilogram in Parisian gravity, and they are then transported to laboratories in other parts of the world, where they are compared with other masses in the local gravity. The quantity defined in this way is called *gravitational mass*.

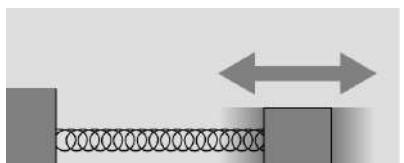
A second and completely different approach is to measure how hard it is to change an object's state of motion. This tells us its *inertial mass*. For example, I'd be more willing to stand in the way of an oncoming poodle than in the path of a freight train, because my body will have a harder time convincing the freight train to stop. This is a dictionary-style conceptual definition, but in physics we need to back up a conceptual definition with an operational definition, which is one that spells out the operations required in order to measure the quantity being defined. We can operationalize our definition of inertial mass by throwing a standard kilogram at an object at a speed of 1 m/s (one meter per second) and measuring the recoiling object's velocity. Suppose we want to measure the mass of a particular block of cement. We put the block in a toy wagon on the sidewalk, and throw a standard kilogram at it. Suppose the standard kilogram hits the wagon, and then drops straight down to the sidewalk, having lost all its velocity, and the wagon and the block inside recoil at a velocity of 0.23 m/s. We then repeat the experiment with the block replaced by various numbers of standard kilograms, and find that we can reproduce the recoil velocity of 0.23 m/s with four standard kilograms in the wagon. We have determined the mass of the block to be four kilograms.¹ Although this definition of inertial mass has an appealing conceptual simplicity, it is obviously not very practical, at least in this crude form. Nevertheless, this method of collision is very much like the methods used for measuring the masses of subatomic particles, which, after all, can't be put on little postal scales!

Astronauts spending long periods of time in space need to monitor their loss of bone and muscle mass, and here as well, it's impossible to measure gravitational mass. Since they don't want to have standard kilograms thrown at them, they use a slightly different technique (figures d and e). They strap themselves to a chair which is attached to a large spring, and measure the time it takes for one cycle of vibration.

¹You might think intuitively that the recoil velocity should be exactly one fourth of a meter per second, and you'd be right except that the wagon has some mass as well. Our present approach, however, only requires that we give a way to test for equality of masses. To predict the recoil velocity from scratch, we'd need to use conservation of momentum, which is discussed in a later chapter.



c / A measurement of inertial mass: the wagon recoils with the same velocity in experiments 1 and 2, establishing that the inertial mass of the cement block is four kilograms.



d / The time for one cycle of vibration is related to the object's inertial mass.



e / Astronaut Tamara Jernigan measures her inertial mass aboard the Space Shuttle.

1.1.1 Problem-solving techniques

How do we use a conservation law, such as conservation of mass, to solve problems? There are two basic techniques.

As an analogy, consider conservation of money, which makes it illegal for you to create dollar bills using your own inkjet printer. (Most people don't intentionally destroy their dollar bills, either!) Suppose the police notice that a particular store doesn't seem to have any customers, but the owner wears lots of gold jewelry and drives a BMW. They suspect that the store is a front for some kind of crime, perhaps counterfeiting. With intensive surveillance, there are two basic approaches they could use in their investigation. One method would be to have undercover agents try to find out how much money goes in the door, and how much money comes back out at the end of the day, perhaps by arranging through some trick to get access to the owner's briefcase in the morning and evening. If the amount of money that comes out every day is greater than the amount that went in, and if they're convinced there is no safe on the premises holding a large reservoir of money, then the owner must be counterfeiting. This inflow-equals-outflow technique is useful if we are sure that there is a region of space within which there is no supply of mass that is being built up or depleted.



f / Example 1.

A stream of water

example 1

If you watch water flowing out of the end of a hose, you'll see that the stream of water is fatter near the mouth of the hose, and skinnier lower down. This is because the water speeds up as it falls. If the cross-sectional area of the stream was equal all along its length, then the rate of flow (kilograms per second) through a lower cross-section would be greater than the rate of flow through a cross-section higher up. Since the flow is steady, the amount of water between the two cross-sections stays constant. Conservation of mass therefore requires that the cross-sectional area of the stream shrink in inverse proportion to the increasing speed of the falling water.

self-check A

Suppose the you point the hose straight up, so that the water is rising rather than falling. What happens as the velocity gets smaller? What happens when the velocity becomes zero? ▷ Answer, p. 1058

How can we apply a conservation law, such as conservation of mass, in a situation where mass might be stored up somewhere? To use a crime analogy again, a prison could contain a certain number of prisoners, who are not allowed to flow in or out at will. In physics, this is known as a *closed system*. A guard might notice that a certain prisoner's cell is empty, but that doesn't mean he's escaped. He could be sick in the infirmary, or hard at work in the shop earning cigarette money. What prisons actually do is to count all their prisoners every day, and make sure today's total is the same as

yesterday's. One way of stating a conservation law is that for a closed system, the total amount of stuff (mass, in this chapter) stays constant.

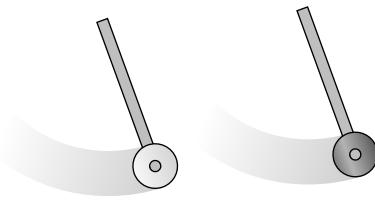
Lavoisier and chemical reactions in a closed system example 2

The French chemist Antoine-Laurent Lavoisier is considered the inventor of the concept of conservation of mass. Before Lavoisier, chemists had never systematically weighed their chemicals to quantify the amount of each substance that was undergoing reactions. They also didn't completely understand that gases were just another state of matter, and hadn't tried performing reactions in sealed chambers to determine whether gases were being consumed from or released into the air. For this they had at least one practical excuse, which is that if you perform a gas-releasing reaction in a sealed chamber with no room for expansion, you get an explosion! Lavoisier invented a balance that was capable of measuring milligram masses, and figured out how to do reactions in an upside-down bowl in a basin of water, so that the gases could expand by pushing out some of the water. In a crucial experiment, Lavoisier heated a red mercury compound, which we would now describe as mercury oxide (HgO), in such a sealed chamber. A gas was produced (Lavoisier later named it "oxygen"), driving out some of the water, and the red compound was transformed into silvery liquid mercury metal. The crucial point was that the total mass of the entire apparatus was exactly the same before and after the reaction. Based on many observations of this type, Lavoisier proposed a general law of nature, that mass is always conserved. (In earlier experiments, in which closed systems were not used, chemists had become convinced that there was a mysterious substance, phlogiston, involved in combustion and oxidation reactions, and that phlogiston's mass could be positive, negative, or zero depending on the situation!)

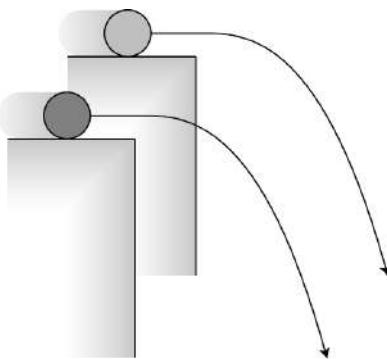
1.1.2 Delta notation

A convenient notation used throughout physics is Δ , the uppercase Greek letter delta, which indicates "change in" or "after minus before." For example, if b represents how much money you have in the bank, then a deposit of \$100 could be represented as $\Delta b = \$100$. That is, the change in your balance was \$100, or the balance after the transaction minus the balance before the transaction equals \$100. A withdrawal would be indicated by $\Delta b < 0$. We represent "before" and "after" using the subscripts i (initial) and f (final), e.g., $\Delta b = b_f - b_i$. Often the delta notation allows more precision than English words. For instance, "time" can be used to mean a point in time ("now's the time"), t , or it could mean a period of time ("the whole time, he had spit on his chin"), Δt .

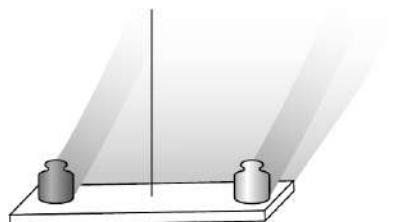
This notation is particularly convenient for discussing conserved quantities. The law of conservation of mass can be stated simply as



a / The two pendulum bobs are constructed with equal gravitational masses. If their inertial masses are also equal, then each pendulum should take exactly the same amount of time per swing.



b / If the cylinders have slightly unequal ratios of inertial to gravitational mass, their trajectories will be a little different.



c / A simplified drawing of an Eötvös-style experiment. If the two masses, made out of two different substances, have slightly different ratios of inertial to gravitational mass, then the apparatus will twist slightly as the earth spins.

$\Delta m = 0$, where m is the total mass of any closed system.

self-check B

If x represents the location of an object moving in one dimension, then how would positive and negative signs of Δx be interpreted? ▷

Answer, p. 1058

Discussion Questions

A If an object had a straight-line $x - t$ graph with $\Delta x = 0$ and $\Delta t \neq 0$, what would be true about its velocity? What would this look like on a graph? What about $\Delta t = 0$ and $\Delta x \neq 0$?

1.2 Equivalence of gravitational and inertial mass

We find experimentally that both gravitational and inertial mass are conserved to a high degree of precision for a great number of processes, including chemical reactions, melting, boiling, soaking up water with a sponge, and rotting of meat and vegetables. Now it's logically possible that both gravitational and inertial mass are conserved, but that there is no particular relationship between them, in which case we would say that they are separately conserved. On the other hand, the two conservation laws may be redundant, like having one law against murder and another law against killing people!

Here's an experiment that gets at the issue: stand up now and drop a coin and one of your shoes side by side. I used a 400-gram shoe and a 2-gram penny, and they hit the floor at the same time as far as I could tell by eye. This is an interesting result, but a physicist and an ordinary person will find it interesting for different reasons.

The layperson is surprised, since it would seem logical that heavier objects would always fall faster than light ones. However, it's fairly easy to prove that if air friction is negligible, any two objects made of the same substance must have identical motion when they fall. For instance, a 2-kg copper mass must exhibit the same falling motion as a 1-kg copper mass, because nothing would be changed by physically joining together two 1-kg copper masses to make a single 2-kg copper mass. Suppose, for example, that they are joined with a dab of glue; the glue isn't under any strain, because the two masses are doing the same thing side by side. Since the glue isn't really doing anything, it makes no difference whether the masses fall separately or side by side.²

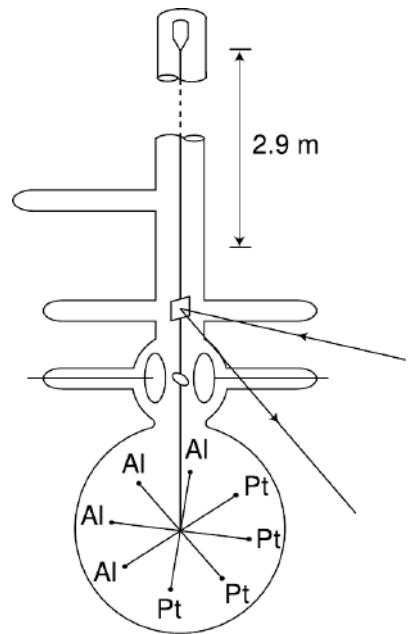
What a physicist finds remarkable about the shoe-and-penny experiment is that it came out the way it did even though the shoe and the penny are made of *different* substances. There is absolutely no theoretical reason why this should be true. We could say that it

²The argument only fails for objects light enough to be affected appreciably by air friction: a bunch of feathers falls differently if you wad them up because the pattern of air flow is altered by putting them together.

happens because the greater gravitational mass of the shoe is exactly counteracted by its greater inertial mass, which makes it harder for gravity to get it moving, but that just leaves us wondering why inertial mass and gravitational mass are always in proportion to each other. It's possible that they are only approximately equivalent. Most of the mass of ordinary matter comes from neutrons and protons, and we could imagine, for instance, that neutrons and protons do not have exactly the same ratio of gravitational to inertial mass. This would show up as a different ratio of gravitational to inertial mass for substances containing different proportions of neutrons and protons.

Galileo did the first numerical experiments on this issue in the seventeenth century by rolling balls down inclined planes, although he didn't think about his results in these terms. A fairly easy way to improve on Galileo's accuracy is to use pendulums with bobs made of different materials. Suppose, for example, that we construct an aluminum bob and a brass bob, and use a double-pan balance to verify to good precision that their gravitational masses are equal. If we then measure the time required for each pendulum to perform a hundred cycles, we can check whether the results are the same. If their inertial masses are unequal, then the one with a smaller inertial mass will go through each cycle faster, since gravity has an easier time accelerating and decelerating it. With this type of experiment, one can easily verify that gravitational and inertial mass are proportional to each other to an accuracy of 10^{-3} or 10^{-4} .

In 1889, the Hungarian physicist Roland Eötvös used a slightly different approach to verify the equivalence of gravitational and inertial mass for various substances to an accuracy of about 10^{-8} , and the best such experiment, figure d, improved on even this phenomenal accuracy, bringing it to the 10^{-12} level.³ In all the experiments described so far, the two objects move along similar trajectories: straight lines in the penny-and-shoe and inclined plane experiments, and circular arcs in the pendulum version. The Eötvös-style experiment looks for differences in the objects' trajectories. The concept can be understood by imagining the following simplified version. Suppose, as in figure b, we roll a brass cylinder off of a tabletop and measure where it hits the floor, and then do the same with an aluminum cylinder, making sure that both of them go over the edge with precisely the same velocity. An object with zero gravitational mass would fly off straight and hit the wall, while an object with zero inertial mass would make a sudden 90-degree turn and drop straight to the floor. If the aluminum and brass cylinders have ordinary, but slightly unequal, ratios of gravitational to inertial mass, then they will follow trajectories that are just slightly different. In other words, if inertial and gravitational mass are not exactly proportional to each other for all substances, then objects made of



d / A more realistic drawing of Braginskii and Panov's experiment. The whole thing was encased in a tall vacuum tube, which was placed in a sealed basement whose temperature was controlled to within 0.02°C . The total mass of the platinum and aluminum test masses, plus the tungsten wire and the balance arms, was only 4.4 g. To detect tiny motions, a laser beam was bounced off of a mirror attached to the wire. There was so little friction that the balance would have taken on the order of several years to calm down completely after being put in place; to stop these vibrations, static electrical forces were applied through the two circular plates to provide very gentle twists on the ellipsoidal mass between them. After Braginskii and Panov.

³V.B. Braginskii and V.I. Panov, Soviet Physics JETP 34, 463 (1972).

different substances will have different trajectories in the presence of gravity.

A simplified drawing of a practical, high-precision experiment is shown in figure c. Two objects made of different substances are balanced on the ends of a bar, which is suspended at the center from a thin fiber. The whole apparatus moves through space on a complicated, looping trajectory arising from the rotation of the earth superimposed on the earth's orbital motion around the sun. Both the earth's gravity and the sun's gravity act on the two objects. If their inertial masses are not exactly in proportion to their gravitational masses, then they will follow slightly different trajectories through space, which will result in a very slight twisting of the fiber between the daytime, when the sun's gravity is pulling upward, and the night, when the sun's gravity is downward. Figure d shows a more realistic picture of the apparatus.

This type of experiment, in which one expects a null result, is a tough way to make a career as a scientist. If your measurement comes out as expected, but with better accuracy than other people had previously achieved, your result is publishable, but won't be considered earthshattering. On the other hand, if you build the most sensitive experiment ever, and the result comes out contrary to expectations, you're in a scary situation. You could be right, and earn a place in history, but if the result turns out to be due to a defect in your experiment, then you've made a fool of yourself.

1.3 Galilean relativity

I defined inertial mass conceptually as a measure of how hard it is to *change* an object's state of motion, the implication being that if you don't interfere, the object's motion won't change. Most people, however, believe that objects in motion have a natural tendency to slow down. Suppose I push my refrigerator to the west for a while at 0.1 m/s, and then stop pushing. The average person would say fridge just naturally stopped moving, but let's imagine how someone in China would describe the fridge experiment carried out in my house here in California. Due to the rotation of the earth, California is moving to the east at about 400 m/s. A point in China at the same latitude has the same speed, but since China is on the other side of the planet, China's east is my west. (If you're finding the three-dimensional visualization difficult, just think of China and California as two freight trains that go past each other, each traveling at 400 m/s.) If I insist on thinking of my dirt as being stationary, then China and its dirt are moving at 800 m/s to my west. From China's point of view, however, it's California that is moving 800 m/s in the opposite direction (my east). When I'm pushing the fridge to the west at 0.1 m/s, the observer in China describes its speed as 799.9 m/s. Once I stop pushing, the fridge speeds back up to 800



a / Galileo Galilei (1564-1642).

m/s. From my point of view, the fridge “naturally” slowed down when I stopped pushing, but according to the observer in China, it “naturally” sped up!

What’s really happening here is that there’s a tendency, due to friction, for the fridge to stop moving *relative to the floor*. In general, only relative motion has physical significance in physics, not absolute motion. It’s not even possible to define absolute motion, since there is no special reference point in the universe that everyone can agree is at rest. Of course if we want to measure motion, we do have to pick some arbitrary reference point which we will say is standing still, and we can then define x , y , and z coordinates extending out from that point, which we can define as having $x = 0$, $y = 0$, $z = 0$. Setting up such a system is known as choosing a *frame of reference*. The local dirt is a natural frame of reference for describing a game of basketball, but if the game was taking place on the deck of a moving ocean liner, we would probably pick a frame of reference in which the deck was at rest, and the land was moving.

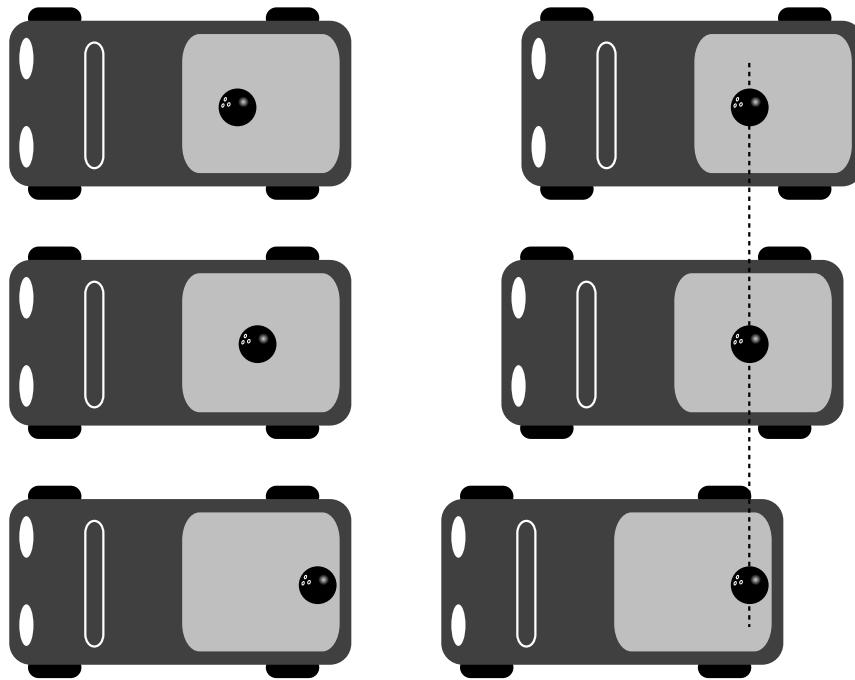
Galileo was the first scientist to reason along these lines, and we now use the term Galilean relativity to refer to a somewhat modernized version of his principle. Roughly speaking, the principle of Galilean relativity states that the same laws of physics apply in any frame of reference that is moving in a straight line at constant speed. We need to refine this statement, however, since it is not necessarily obvious which frames of reference are going in a straight line at constant speed. A person in a pickup truck pulling away from a stoplight could admit that the car’s velocity is changing, or she could insist that the truck is at rest, and the meter on the dashboard is going up because the asphalt picked that moment to start moving faster and faster backward! Frames of reference are not all created equal, however, and the accelerating truck’s frame of reference is not as good as the asphalt’s. We can tell, because a bowling ball in the back of the truck, as in figure c, appears to behave strangely in the driver’s frame of reference: in her rear-view mirror, she sees the ball, initially at rest, start moving faster and faster toward the back of the truck. This goofy behavior is evidence that there is something wrong with her frame of reference. A person on the sidewalk, however, sees the ball as standing still. In the sidewalk’s frame of reference, the truck pulls away from the ball, and this makes sense, because the truck is burning gas and using up energy to change its state of motion.

We therefore define an *inertial frame of reference* as one in which we never see objects change their state of motion without any apparent reason. The sidewalk is a pretty good inertial frame, and a car moving relative to the sidewalk at constant speed in a straight line defines a pretty good inertial frame, but a car that is accelerating or turning is not a inertial frame.



b / The earth spins. People in Shanghai say they’re at rest and people in Los Angeles are moving. Angelenos say the same about the Shanghaiese.

c / Left: In a frame of reference that speeds up with the truck, the bowling ball appears to change its state of motion for no reason. Right: In an inertial frame of reference, which the surface of the earth approximately is, the bowling ball stands still, which makes sense because there is nothing that would cause it to change its state of motion.



The principle of Galilean relativity states that inertial frames exist, and that the same laws of physics apply in all inertial frames of reference, regardless of one frame's straight-line, constant-speed motion relative to another.⁴

Another way of putting it is that all inertial frames are created equal. We can say whether one inertial frame is in motion or at rest relative to another, but there is no privileged “rest frame.” There is no experiment that comes out any different in laboratories in different inertial frames, so there is no experiment that could tell us which inertial frame is really, truly at rest.

The speed of sound

example 3

- ▷ The speed of sound in air is only 340 m/s, so unless you live at a near-polar latitude, you’re moving at greater than the speed of sound right now due to the Earth’s rotation. In that case, why don’t we experience exciting phenomena like sonic booms all the time? ▷ It might seem as though you’re unprepared to deal with this question right now, since the only law of physics you know is conservation of mass, and conservation of mass doesn’t tell you anything obviously useful about the speed of sound or sonic booms. Galilean relativity, however, is a blanket statement about all the laws of physics, so in a situation like this, it may let you predict the results of the laws of physics without actually knowing what all the laws are! If the laws of physics predict a certain value for the speed of sound, then they had better predict the speed

⁴The principle of Galilean relativity is extended on page 195.

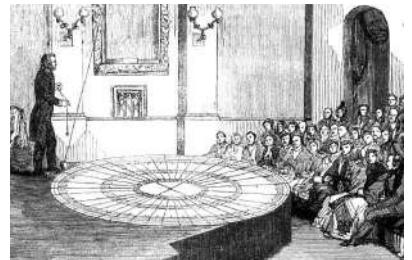
of the sound relative to the air, not their speed relative to some special “rest frame.” Since the air is moving along with the rotation of the earth, we don’t detect any special phenomena. To get a sonic boom, the source of the sound would have to be moving relative to the air.

The Foucault pendulum

example 4

Note that in the example of the bowling ball in the truck, I didn’t claim the sidewalk was *exactly* a Galilean frame of reference. This is because the sidewalk is moving in a circle due to the rotation of the Earth, and is therefore changing the direction of its motion continuously on a 24-hour cycle. However, the curve of the motion is so gentle that under ordinary conditions we don’t notice that the local dirt’s frame of reference isn’t quite inertial. The first demonstration of the noninertial nature of the earth-fixed frame of reference was by Foucault using a very massive pendulum (figure d) whose oscillations would persist for many hours without becoming imperceptible. Although Foucault did his demonstration in Paris, it’s easier to imagine what would happen at the north pole: the pendulum would keep swinging in the same plane, but the earth would spin underneath it once every 24 hours. To someone standing in the snow, it would appear that the pendulum’s plane of motion was twisting. The effect at latitudes less than 90 degrees turns out to be slower, but otherwise similar. The Foucault pendulum was the first definitive experimental proof that the earth really did spin on its axis, although scientists had been convinced of its rotation for a century based on more indirect evidence about the structure of the solar system.

Although popular belief has Galileo being prosecuted by the Catholic Church for saying the earth rotated on its axis and also orbited the sun, Foucault’s pendulum was still centuries in the future, so Galileo had no hard proof; Galileo’s insights into relative versus absolute motion simply made it more plausible that the world could be spinning without producing dramatic effects, but didn’t disprove the contrary hypothesis that the sun, moon, and stars went around the earth every 24 hours. Furthermore, the Church was much more liberal and enlightened than most people believe. It didn’t (and still doesn’t) require a literal interpretation of the Bible, and one of the Church officials involved in the Galileo affair wrote that “the Bible tells us how to go to heaven, not how the heavens go.” In other words, religion and science should be separate. The actual reason Galileo got in trouble is shrouded in mystery, since Italy in the age of the Medicis was a secretive place where unscrupulous people might settle a score with poison or a false accusation of heresy. What is certain is that Galileo’s satirical style of scientific writing made many enemies among the powerful Jesuit scholars who were his intellectual opponents — he compared one to a snake that doesn’t know its own back is broken. It’s also possible that the Church was far

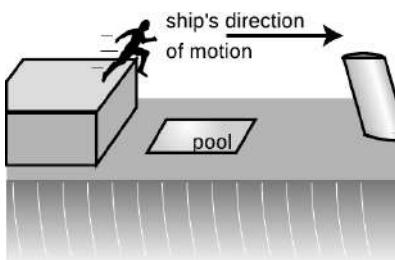


d / Foucault demonstrates his pendulum to an audience at a lecture in 1851.



e / Galileo’s trial.

less upset by his astronomical work than by his support for atomism (discussed further in the next section). Some theologians perceived atomism as contradicting transubstantiation, the Church's doctrine that the holy bread and wine were literally transformed into the flesh and blood of Christ by the priest's blessing.



f / Discussion question A.

self-check C

What is incorrect about the following supposed counterexamples to the principle of inertia?

(1) When astronauts blast off in a rocket, their huge velocity does cause a physical effect on their bodies — they get pressed back into their seats, the flesh on their faces gets distorted, and they have a hard time lifting their arms.

(2) When you're driving in a convertible with the top down, the wind in your face is an observable physical effect of your absolute motion. ▷

Answer, p. 1058

▷ *Solved problem: a bug on a wheel*

page 71, problem 12

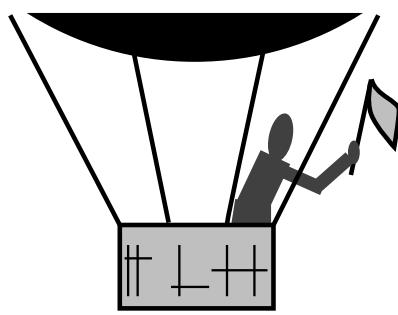
Discussion Questions

A A passenger on a cruise ship finds, while the ship is docked, that he can leap off of the upper deck and just barely make it into the pool on the lower deck. If the ship leaves dock and is cruising rapidly, will this adrenaline junkie still be able to make it?

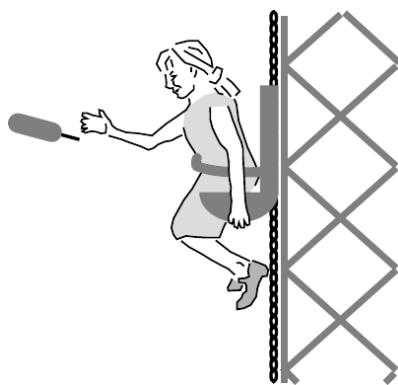
B You are a passenger in the open basket hanging under a helium balloon. The balloon is being carried along by the wind at a constant velocity. If you are holding a flag in your hand, will the flag wave? If so, which way? [Based on a question from PSSC Physics.]

C Aristotle stated that all objects naturally wanted to come to rest, with the unspoken implication that "rest" would be interpreted relative to the surface of the earth. Suppose we could transport Aristotle to the moon, put him in a space suit, and kick him out the door of the spaceship and into the lunar landscape. What would he expect his fate to be in this situation? If intelligent creatures inhabited the moon, and one of them independently came up with the equivalent of Aristotelian physics, what would they think about objects coming to rest?

D Sally is on an amusement park ride which begins with her chair being hoisted straight up a tower at a constant speed of 60 miles/hour. Despite stern warnings from her father that he'll take her home the next time she misbehaves, she decides that as a scientific experiment she really needs to release her corndog over the side as she's on the way up. She does not throw it. She simply sticks it out of the car, lets it go, and watches it against the background of the sky, with no trees or buildings as reference points. What does the corndog's motion look like as observed by Sally? Does its speed ever appear to her to be zero? What acceleration does she observe it to have: is it ever positive? negative? zero? What would her enraged father answer if asked for a similar description of its motion as it appears to him, standing on the ground?



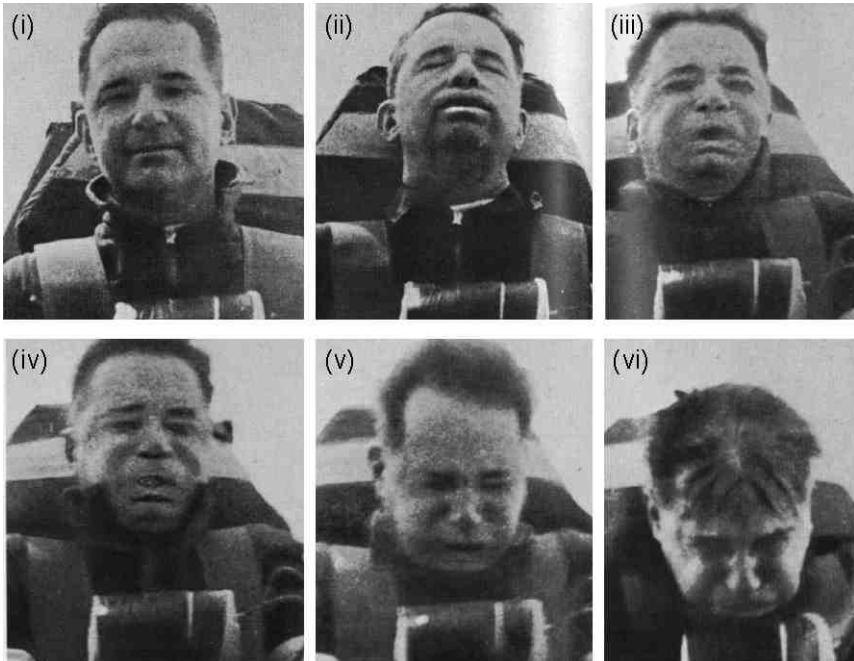
g / Discussion question B.



h / Discussion question D.

1.3.1 Applications of calculus

Let's see how this relates to calculus. If an object is moving in one dimension, we can describe its position with a function $x(t)$. The derivative $v = dx/dt$ is called the velocity, and the second derivative $a = dv/dt = d^2x/dt^2$ is the acceleration. Galilean relativity tells us that there is no detectable effect due to an object's absolute velocity, since in some other frame of reference, the object's velocity might be zero. However, an acceleration does have physical consequences.

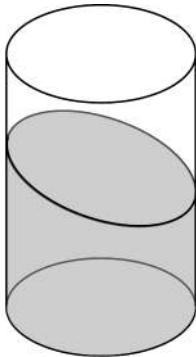


i / This Air Force doctor volunteered to ride a rocket sled as a medical experiment. The obvious effects on his head and face are not because of the sled's speed but because of its rapid *changes* in speed: increasing in (ii) and (iii), and decreasing in (v) and (vi). In (iv) his speed is greatest, but because his speed is not increasing or decreasing very much at this moment, there is little effect on him. (U.S. Air Force)

Observers in different inertial frames of reference will disagree on velocities, but agree on accelerations. Let's keep it simple by continuing to work in one dimension. One frame of reference uses a coordinate system x_1 , and the other we label x_2 . If the positive x_1 and x_2 axes point in the same direction, then in general two inertial frames could be related by an equation of the form $x_2 = x_1 + b + ut$, where u is the constant velocity of one frame relative to the other, and the constant b tells us how far apart the origins of the two coordinate systems were at $t = 0$. The velocities are different in the two frames of reference:

$$\frac{dx_2}{dt} = \frac{dx_1}{dt} + u,$$

Suppose, for example, frame 1 is defined from the sidewalk, and frame 2 is fixed to a float in a parade that is moving to our left at a velocity $u = 1 \text{ m/s}$. A dog that is moving to the right with a velocity $v_1 = dx_1/dt = 3 \text{ m/s}$ in the sidewalk's frame will appear to be moving at a velocity of $v_2 = dx_2/dt = dx_1/dt + u = 4 \text{ m/s}$ in the float's frame.



Self-check D.

For acceleration, however, we have

$$\frac{d^2 x_2}{dt^2} = \frac{d^2 x_1}{dt^2},$$

since the derivative of the constant u is zero. Thus an acceleration, unlike a velocity, can have a definite physical significance to all observers in all frames of reference. If this wasn't true, then there would be no particular reason to define a quantity called acceleration in the first place.

self-check D

The figure shows a bottle of beer sitting on a table in the dining car of a train. Does the tilting of the surface tell us about the train's velocity, or its acceleration? What would a person in the train say about the bottle's velocity? What about a person standing in a field outside and looking in through the window? What about the acceleration?

▷ Answer, p. 1058

1.4 A preview of some modern physics

"Mommy, why do you and Daddy have to go to work?" "To make money, sweetie-pie." "Why do we need money?" "To buy food." "Why does food cost money?" When small children ask a chain of "why" questions like this, it usually isn't too long before their parents end up saying something like, "Because that's just the way it is," or, more honestly, "I don't know the answer."

The same happens in physics. We may gradually learn to explain things more and more deeply, but there's always the possibility that a certain observed fact, such as conservation of mass, will never be understood on any deeper level. Science, after all, uses limited methods to achieve limited goals, so the ultimate reason for all existence will always be the province of religion. There is, however, an appealing explanation for conservation of mass, which is atomism, the theory that matter is made of tiny, unchanging particles. The atomic hypothesis dates back to ancient Greece, but the first solid evidence to support it didn't come until around the eighteenth century, and individual atoms were never detected until about 1900. The atomic theory implies not only conservation of mass, but a couple of other things as well.

First, it implies that the total mass of one particular element is conserved. For instance, lead and gold are both elements, and if we assume that lead atoms can't be turned into gold atoms, then the total mass of lead and the total mass of gold are separately conserved. It's as though there was not just a law against pickpocketing, but also a law against surreptitiously moving money from one of the victim's pockets to the other. It turns out, however, that although chemical reactions never change one type of atom into another, transmutation can happen in nuclear reactions, such as the

ones that created most of the elements in your body out of the primordial hydrogen and helium that condensed out of the aftermath of the Big Bang.

Second, atomism implies that mass is *quantized*, meaning that only certain values of mass are possible and the ones in between can't exist. We can have three atoms of gold or four atoms of gold, but not three and a half. Although quantization of mass is a natural consequence of any theory in which matter is made up of tiny particles, it was discovered in the twentieth century that other quantities, such as energy, are quantized as well, which had previously not been suspected.

self-check E

Is money quantized?

► Answer, p. 1058

If atomism is starting to make conservation of mass seem inevitable to you, then it may disturb you to know that Einstein discovered it isn't really conserved. If you put a 50-gram iron nail in some water, seal the whole thing up, and let it sit on a fantastically precise balance while the nail rusts, you'll find that the system loses about 6×10^{-12} kg of mass by the time the nail has turned completely to rust. This has to do with Einstein's famous equation $E = mc^2$. Rusting releases heat energy, which then escapes out into the room. Einstein's equation states that this amount of heat, E , is equivalent to a certain amount of mass, m . The c in the c^2 is the speed of light, which is a large number, and a large amount of energy is therefore equivalent to a very small amount of mass, so you don't notice nonconservation of mass under ordinary conditions. What is really conserved is not the mass, m , but the mass-plus-energy, $E + mc^2$. The point of this discussion is not to get you to do numerical exercises with $E = mc^2$ (at this point you don't even know what units are used to measure energy), but simply to point out to you the empirical nature of the laws of physics. If a previously accepted theory is contradicted by an experiment, then the theory needs to be changed. This is also a good example of something called the *correspondence principle*, which is a historical observation about how scientific theories change: when a new scientific theory replaces an old one, the old theory is always contained within the new one as an approximation that works within a certain restricted range of situations. Conservation of mass is an extremely good approximation for all chemical reactions, since chemical reactions never release or consume enough energy to change the total mass by a large percentage. Conservation of mass would not have been accepted for 110 years as a fundamental principle of physics if it hadn't been verified over and over again by a huge number of accurate experiments.

This chapter is summarized on page 1075. Notation and terminology are tabulated on pages 1070-1071.

Problems

The symbols \checkmark , \blacksquare , etc. are explained on page 72.

1 Thermometers normally use either mercury or alcohol as their working fluid. If the level of the fluid rises or falls, does this violate conservation of mass? \blacksquare

2 The ratios of the masses of different types of atoms were determined a century before anyone knew any actual atomic masses in units of kg. One finds, for example, that when ordinary table salt, NaCl, is melted, the chlorine atoms bubble off as a gas, leaving liquid sodium metal. Suppose the chlorine escapes, so that its mass cannot be directly determined by weighing. Experiments show that when 1.00000 kg of NaCl is treated in this way, the mass of the remaining sodium metal is 0.39337 kg. Based on this information, determine the ratio of the mass of a chlorine atom to that of a sodium atom. \checkmark \blacksquare

3 An atom of the most common naturally occurring uranium isotope breaks up spontaneously into a thorium atom plus a helium atom. The masses are as follows:

uranium	$3.95292849 \times 10^{-25}$ kg
thorium	$3.88638748 \times 10^{-25}$ kg
helium	6.646481×10^{-27} kg

Each of these experimentally determined masses is uncertain in its last decimal place. Is mass conserved in this process to within the accuracy of the experimental data? How would you interpret this? \blacksquare

4 If two spherical water droplets of radius b combine to make a single droplet, what is its radius? (Assume that water has constant density.) \blacksquare

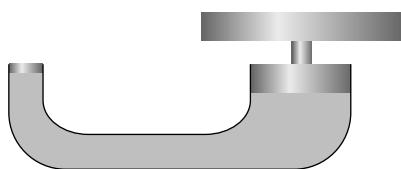
5 Make up an experiment that would test whether mass is conserved in an animal's metabolic processes. \blacksquare

6 The figure shows a hydraulic jack. What is the relationship between the distance traveled by the plunger and the distance traveled by the object being lifted, in terms of the cross-sectional areas? \blacksquare

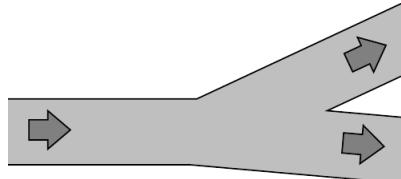
7 In an example in this chapter, I argued that a stream of water must change its cross-sectional area as it rises or falls. Suppose that the stream of water is confined to a constant-diameter pipe. Which assumption breaks down in this situation? \blacksquare

8 A river with a certain width and depth splits into two parts, each of which has the same width and depth as the original river. What can you say about the speed of the current after the split? \blacksquare

9 The diagram shows a cross-section of a wind tunnel of the

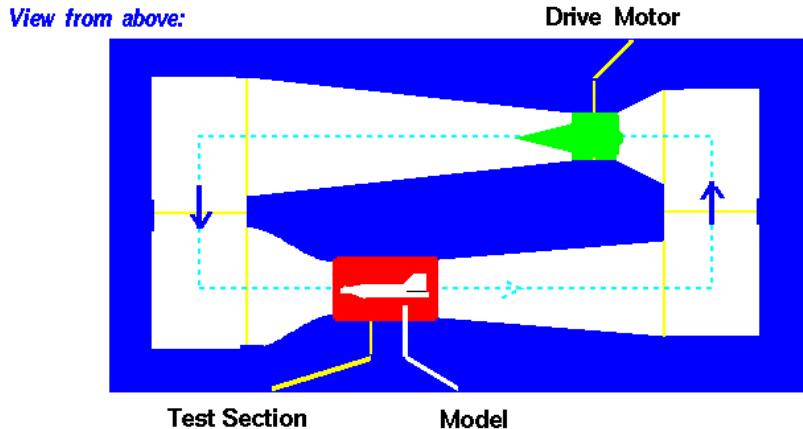


Problem 6.



Problem 8.

kind used, for example, to test designs of airplanes. Under normal conditions of use, the density of the air remains nearly constant throughout the whole wind tunnel. How can the speed of the air be controlled and calculated? (Diagram by NASA, Glenn Research Center.)



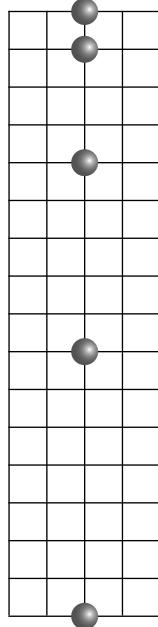
- 10** A water wave is in a tank that extends horizontally from $x = 0$ to $x = a$, and from $z = 0$ to $z = b$. We assume for simplicity that at a certain moment in time the height y of the water's surface only depends on x , not z , so that we can effectively ignore the z coordinate. Under these assumptions, the total volume of the water in the tank is

$$V = b \int_0^a y(x) dx.$$

Since the density of the water is essentially constant, conservation of mass requires that V is always the same. When the water is calm, we have $y = h$, where $h = V/ab$. If two different wave patterns move into each other, we might imagine that they would add in the sense that $y_{total} - h = (y_1 - h) + (y_2 - h)$. Show that this type of addition is consistent with conservation of mass. ■

- 11** The figure shows the position of a falling ball at equal time intervals, depicted in a certain frame of reference. On a similar grid, show how the ball's motion would appear in a frame of reference that was moving horizontally at a speed of one box per unit time relative to the first frame. ■

- 12** The figure shows the motion of a point on the rim of a rolling wheel. (The shape is called a cycloid.) Suppose bug A is riding on the rim of the wheel on a bicycle that is rolling, while bug B is on the spinning wheel of a bike that is sitting upside down on the floor. Bug A is moving along a cycloid, while bug B is moving in a circle. Both wheels are doing the same number of revolutions per minute. Which bug has a harder time holding on, or do they find it equally



Problem 11.



Problem 12.

difficult?

▷ Solution, p. 1039 ■

Key to symbols:

■ easy ■ typical ■ challenging ■ difficult ■ very difficult
✓ An answer check is available at www.lightandmatter.com.

Chapter 2

Conservation of Energy

Do you pronounce it Joule's to rhyme with schools,
Joule's to rhyme with Bowls,
or Joule's to rhyme with Scowls?
Whatever you call it, by Joule's,
or Joule's,
or Joule's, it's good!

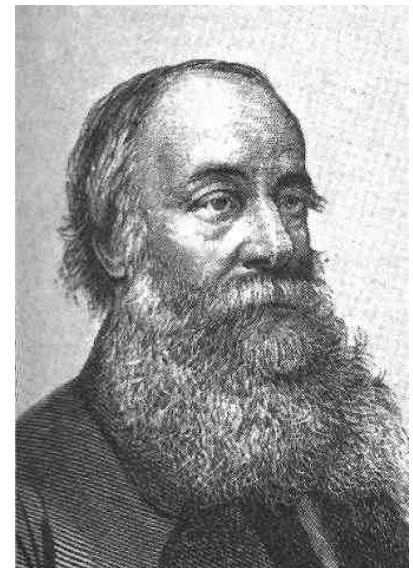
Advertising slogan of the Joule brewery. The name, and the corresponding unit of energy, are now usually pronounced so as to rhyme with “school.”

2.1 Energy

2.1.1 The energy concept

You'd probably like to be able to drive your car and light your apartment without having to pay money for gas and electricity, and if you do a little websurfing, you can easily find people who say they have the solution to your problem. This kind of scam has been around for centuries. It used to be known as a perpetual motion machine, but nowadays the con artists' preferred phrase is "free energy."¹ A typical "free-energy" machine would be a sealed box that heats your house without needing to be plugged into a wall socket or a gas pipe. Heat comes out, but nothing goes in, and this can go on indefinitely. But an interesting thing happens if you try to check on the advertised performance of the machine. Typically, you'll find out that either the device is still in development, or it's back-ordered because so many people have already taken advantage of this Fantastic Opportunity! In a few cases, the magic box exists, but the inventor is only willing to demonstrate very small levels of heat output for short periods of time, in which case there's probably a tiny hearing-aid battery hidden in there somewhere, or some other trick.

Since nobody has ever succeeded in building a device that creates heat out of nothing, we might also wonder whether any device exists



a / James Joule, 1818-1889. The son of a wealthy brewer, Joule was tutored as a young man by the famous scientist John Dalton. Fascinated by electricity, he and his brother experimented by giving electric shocks to each other and to the family's servants. Joule ran the brewery as an adult, and science was merely a serious hobby. His work on energy can be traced to his attempt to build an electric motor that would replace steam engines. His ideas were not accepted at first, partly because they contradicted the widespread belief that heat was a fluid, and partly because they depended on extremely precise measurements, which had not previously been common in physics.

¹An entertaining account of this form of quackery is given in **Voodoo Science: The Road from Foolishness to Fraud**, Robert Park, Oxford University Press, 2000. Until reading this book, I hadn't realized the degree to which pseudoscience had penetrated otherwise respectable scientific organizations like NASA.



b / Heat energy can be converted to light energy. Very hot objects glow visibly, and even objects that aren't so hot give off infrared light, a color of light that lies beyond the red end of the visible rainbow. This photo was made with a special camera that records infrared light. The man's warm skin emits quite a bit of infrared light energy, while his hair, at a lower temperature, emits less.

that can do the opposite, turning heat into nothing. You might think that a refrigerator was such a device, but actually your refrigerator doesn't destroy the heat in the food. What it really does is to extract some of the heat and bring it out into the room. That's why it has big radiator coils on the back, which get hot when it's in operation.

If it's not possible to destroy or create heat outright, then you might start to suspect that heat was a conserved quantity. This would be a successful rule for explaining certain processes, such as the transfer of heat between a cold Martini and a room-temperature olive: if the olive loses a little heat, then the drink must gain the same amount. It would fail in general, however. Sunlight can heat your skin, for example, and a hot lightbulb filament can cool off by emitting light. Based on these observations, we could revise our proposed conservation law, and say that there is something called heatpluslight, which is conserved. Even this, however, needs to be generalized in order to explain why you can get a painful burn playing baseball when you slide into a base. Now we could call it heatpluslightplusmotion. The word is getting pretty long, and we haven't even finished the list.

Rather than making the word longer and longer, physicists have hijacked the word "energy" from ordinary usage, and give it a new, specific technical meaning. Just as the Parisian platinum-iridium kilogram defines a specific unit of mass, we need to pick something that defines a definite unit of energy. The metric unit of energy is the joule (J), and we'll define it as the amount of energy required to heat 0.24 grams of water from 20 to 21 degrees Celsius. (Don't memorize the numbers.)²

Temperature of a mixture

example 1

- ▷ If 1.0 kg of water at 20°C is mixed with 4.0 kg of water at 30°C, what is the temperature of the mixture?

- ▷ Let's assume as an approximation that each degree of temperature change corresponds to the same amount of energy. In other words, we assume $\Delta E = mc\Delta T$, regardless of whether, as in the definition of the joule, we have $\Delta T = 21^\circ\text{C}-20^\circ\text{C}$ or, as in the present example, some other combination of initial and final temperatures. To be consistent with the definition of the joule, we must have $c = (1 \text{ J})/(0.24 \text{ g})/(1^\circ\text{C}) = 4.2 \times 10^3 \text{ J/kg}\cdot^\circ\text{C}$, which is referred to as the specific heat of water.

²Although the definition refers to the Celsius scale of temperature, it's not necessary to give an operational definition of the temperature concept in general (which turns out to be quite a tricky thing to do completely rigorously); we only need to establish two specific temperatures that can be reproduced on thermometers that have been calibrated in a standard way. Heat and temperature are discussed in more detail in section 2.4, and in chapter 5. Conceptually, heat is a measure of energy, whereas temperature relates to how concentrated that energy is.

Conservation of energy tells us $\Delta E = 0$, so

$$m_1 c \Delta T_1 + m_2 c \Delta T_2 = 0$$

or

$$\begin{aligned}\frac{\Delta T_1}{\Delta T_2} &= -\frac{m_2}{m_1} \\ &= -4.0\end{aligned}$$

If T_1 has to change four times as much as T_2 , and the two final temperatures are equal, then the final temperature must be 28°C .

Note how only *differences* in temperature and energy appeared in the preceding example. In other words, we don't have to make any assumptions about whether there is a temperature at which all an object's heat energy is removed. Historically, the energy and temperature units were invented before it was shown that there is such a temperature, called absolute zero. There is a scale of temperature, the Kelvin scale, in which the unit of temperature is the same as the Celsius degree, but the zero point is defined as absolute zero. But as long as we only deal with temperature differences, it doesn't matter whether we use Kelvin or Celsius. Likewise, as long as we deal with differences in heat energy, we don't normally have to worry about the total amount of heat energy the object has. In standard physics terminology, "heat" is used only to refer to differences, while the total amount is called the object's "thermal energy." This distinction is often ignored by scientists in casual speech, and in this book I'll usually use "heat" for either quantity.

We're defining energy by adding up things from a list, which we lengthen as needed: heat, light, motion, etc. One objection to this approach is aesthetic: physicists tend to regard complication as a synonym for ugliness. If we have to keep on adding more and more forms of energy to our laundry list, then it's starting to sound like energy is distressingly complicated. Luckily it turns out that energy is simpler than it seems. Many forms of energy that are apparently unrelated turn out to be manifestations of a small number of forms at the atomic level, and this is the topic of section 2.4.

Discussion Questions

A The ancient Greek philosopher Aristotle said that objects "naturally" tended to slow down, unless there was something pushing on them to keep them moving. What important insight was he missing?

2.1.2 Logical issues

Another possible objection is that the open-ended approach to defining energy might seem like a kind of cheat, since we keep on inventing new forms whenever we need them. If a certain experiment seems to violate conservation of energy, can't we just invent

a new form of invisible “mystery energy” that patches things up? This would be like balancing your checkbook by putting in a fake transaction that makes your calculation of the balance agree with your bank’s. If we could fudge this way, then conservation of energy would be untestable — impossible to prove or disprove.

Actually all scientific theories are unprovable. A theory can never be proved, because the experiments can only cover a finite number out of the infinitely many situations in which the theory is supposed to apply. Even a million experiments won’t suffice to prove it in the same sense of the word “proof” that is used in mathematics. However, even one experiment that contradicts a theory is sufficient to show that the theory is wrong. A theory that is immune to disproof is a bad theory, because there is no way to test it. For instance, if I say that 23 is the maximum number of angels that can dance on the head of a pin, I haven’t made a properly falsifiable scientific theory, since there’s no method by which anyone could even attempt to prove me wrong based on observations or experiments.

Conservation of energy is testable because new forms of energy are expected to show regular mathematical behavior, and are supposed to be related in a measurable way to observable phenomena. As an example, let’s see how to extend the energy concept to include motion.

2.1.3 Kinetic energy

Energy of motion is called *kinetic energy*. (The root of the word is the same as the word “cinema” – in French, kinetic energy is “énergie cinétique.”) How does an object’s kinetic energy depend on its mass and velocity? Joule attempted a conceptually simple experiment on his honeymoon in the French-Swiss Alps near Mt. Chamonix, in which he measured the difference in temperature between the top and bottom of a waterfall. The water at the top of the falls has some gravitational energy, which isn’t our subject right now, but as it drops, that gravitational energy is converted into kinetic energy, and then into heat energy due to internal friction in the churning pool at the bottom:

$$\text{gravitational energy} \rightarrow \text{kinetic energy} \rightarrow \text{heat energy}$$

In the logical framework of this book’s presentation of energy, the significance of the experiment is that it provides a way to find out how an object’s kinetic energy depends on its mass and velocity. The increase in heat energy should equal the kinetic energy of the water just before impact, so in principle we could measure the water’s mass, velocity, and kinetic energy, and see how they relate to one another.³

³From Joule’s point of view, the point of the experiment was different. At that time, most physicists believed that heat was a quantity that was conserved



c / As in figure b, an infrared camera distinguishes hot and cold areas. As the bike skids to a stop with its brakes locked, the kinetic energy of the bike and rider is converted into heat in both the floor (top) and the tire (bottom).

Although the story is picturesque and memorable, most books that mention the experiment fail to note that it was a failure! The problem was that heat wasn't the only form of energy being released. In reality, the situation was more like this:

$$\begin{aligned} \text{gravitational energy} &\rightarrow \text{kinetic energy} \\ &\rightarrow \text{heat energy} \\ &+ \text{sound energy} \\ &+ \text{energy of partial evaporation.} \end{aligned}$$

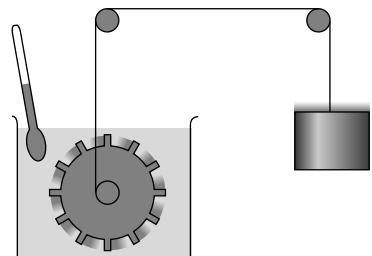
The successful version of the experiment, shown in figures d and f, used a paddlewheel spun by a dropping weight. As with the waterfall experiment, this one involves several types of energy, but the difference is that in this case, they can all be determined and taken into account. (Joule even took the precaution of putting a screen between himself and the can of water, so that the infrared light emitted by his warm body wouldn't warm it up at all!) The result⁴ is

$$K = \frac{1}{2}mv^2 \quad [\text{kinetic energy}].$$

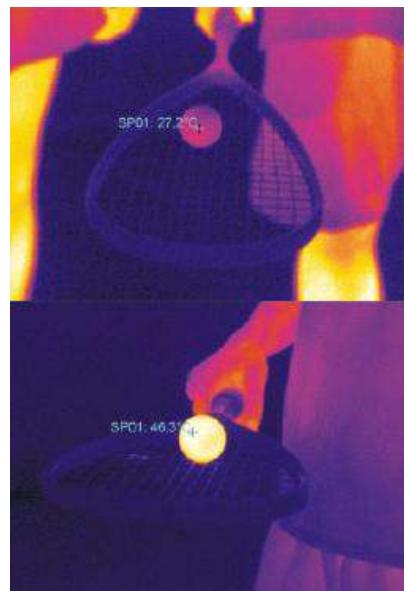
Whenever you encounter an equation like this for the first time, you should get in the habit of interpreting it. First off, we can tell that by making the mass or velocity greater, we'd get more kinetic energy. That makes sense. Notice, however, that we have mass to the first power, but velocity to the second. Having the whole thing proportional to mass to the first power is necessary on theoretical grounds, since energy is supposed to be additive. The dependence on v^2 couldn't have been predicted, but it is sensible. For instance, suppose we reverse the direction of motion. This would reverse the sign of v , because in one dimension we use positive and negative signs to indicate the direction of motion. But since v^2 is what appears in the equation, the resulting kinetic energy is unchanged.

separately from the rest of the things to which we now refer as energy, i.e., mechanical energy. Separate units of measurement had been constructed for heat and mechanical energy, but Joule was trying to show that one could convert back and forth between them, and that it was actually their sum that was conserved, if they were both expressed in consistent units. His main result was the conversion factor that would allow the two sets of units to be reconciled. By showing that the conversion factor came out the same in different types of experiments, he was supporting his assertion that heat was not separately conserved. From Joule's perspective or from ours, the result is to connect the mysterious, invisible phenomenon of heat with forms of energy that are visible properties of objects, i.e., mechanical energy.

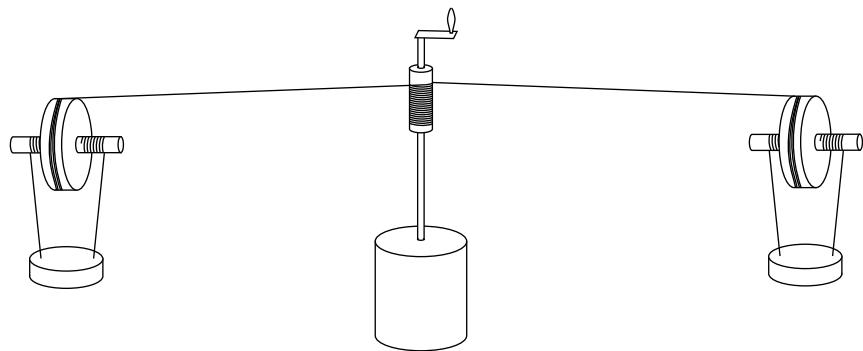
⁴If you've had a previous course in physics, you may have seen this presented not as an empirical result but as a theoretical one, derived from Newton's laws, and in that case you might feel you're being cheated here. However, I'm going to reverse that reasoning and derive Newton's laws from the conservation laws in chapter 3. From the modern perspective, conservation laws are more fundamental, because they apply in cases where Newton's laws don't.



d / A simplified drawing of Joule's paddlewheel experiment.



e / The heating of the tire and floor in figure c is something that the average person might have predicted in advance, but there are other situations where it's not so obvious. When a ball slams into a wall, it doesn't rebound with the same amount of kinetic energy. Was some energy destroyed? No. The ball and the wall heat up. These infrared photos show a squash ball at room temperature (top), and after it has been played with for several minutes (bottom), causing it to heat up detectably.



f / A realistic drawing of Joule's apparatus, based on the illustration in his original paper. The paddlewheel is sealed inside the can in the middle. Joule wound up the two 13-kg lead weights and dropped them 1.6 meters, repeating this 20 times to produce a temperature change of only about half a degree Fahrenheit in the water inside the sealed can. He claimed in his paper to be able to measure temperatures to an accuracy of 1/200 of a degree.

What about the factor of 1/2 in front? It comes out to be exactly 1/2 by the design of the metric system. If we'd been using the old-fashioned British engineering system of units (which is no longer used in the U.K.), the equation would have been $K = (7.44 \times 10^{-2} \text{ Btu} \cdot \text{s}^2/\text{slug} \cdot \text{ft}^2)mv^2$. The version of the metric system called the SI,⁵ in which everything is based on units of kilograms, meters, and seconds, not only has the numerical constant equal to 1/2, but makes it unitless as well. In other words, we can think of the joule as simply an abbreviation, 1 J=1 kg·m²/s². More familiar examples of this type of abbreviation are 1 minute=60 s, and the metric unit of land area, 1 hectare=10000 m².

Ergs and joules

example 2

▷ There used to be two commonly used systems of metric units, referred to as mks and cgs. The mks system, now called the SI, is based on the meter, the kilogram, and the second. The cgs system, which is now obsolete, was based on the centimeter, the gram, and the second. In the cgs system, the unit of energy is not the joule but the erg, 1 erg=1 g · cm²/s². How many ergs are in one joule?

▷ The simplest approach is to treat the units as if they were alge-

⁵Système International

bra symbols.

$$\begin{aligned}
 1 \text{ J} &= 1 \frac{\text{kg} \cdot \text{m}^2}{\text{s}^2} \\
 &= 1 \frac{\text{kg} \cdot \text{m}^2}{\text{s}^2} \times \frac{1000 \text{ g}}{1 \text{ kg}} \times \left(\frac{100 \text{ cm}}{1 \text{ m}} \right)^2 \\
 &= 10^7 \frac{\text{g} \cdot \text{cm}^2}{\text{s}^2} \\
 &= 10^7 \text{ erg}
 \end{aligned}$$

Cabin air in a jet airplane

example 3

▷ A jet airplane typically cruises at a velocity of 270 m/s. Outside air is continuously pumped into the cabin, but must be cooled off first, both because (1) it heats up due to friction as it enters the engines, and (2) it is heated as a side-effect of being compressed to cabin pressure. Calculate the increase in temperature due to the first effect. The specific heat of dry air is about $1.0 \times 10^3 \text{ J/kg} \cdot ^\circ\text{C}$.

▷ This is easiest to understand in the frame of reference of the plane, in which the air rushing into the engine is stopped, and its kinetic energy converted into heat.⁶ Conservation of energy tells us

$$\begin{aligned}
 0 &= \Delta E \\
 &= \Delta K + \Delta E_{heat}.
 \end{aligned}$$

In the plane's frame of reference, the air's initial velocity is $v_i = 270 \text{ m/s}$, and its final velocity is zero, so the change in its kinetic energy is negative,

$$\begin{aligned}
 \Delta K &= K_f - K_i \\
 &= 0 - (1/2)mv_i^2 \\
 &= -(1/2)mv_i^2.
 \end{aligned}$$

Assuming that the specific heat of air is roughly independent of temperature (which is why the number was stated with the word "about"), we can substitute into $0 = \Delta K + \Delta E_{heat}$, giving

$$\begin{aligned}
 0 &= -\frac{1}{2}mv_i^2 + mc\Delta T \\
 \frac{1}{2}v_i^2 &= c\Delta T.
 \end{aligned}$$

Note how the mass cancels out. This is a big advantage of solving problems algebraically first, and waiting until the end to plug in

⁶It's not at all obvious that the solution would work out in the earth's frame of reference, although Galilean relativity states that it doesn't matter which frame we use. Chapter 3 discusses the relationship between conservation of energy and Galilean relativity.

numbers. With a purely numerical approach, we wouldn't even have known what value of m to pick, or if we'd guessed a value like 1 kg, we wouldn't have known whether our answer depended on that guess.

Solving for ΔT , and writing v instead of v_i for simplicity, we find

$$\Delta T = \frac{v^2}{2c}$$

$$\approx 40^\circ\text{C}.$$

The passengers would be boiled alive if not for the refrigeration. The first stage of cooling happens via heat exchangers in the engine struts, but a second stage, using a refrigerator under the floor of the cabin, is also necessary. Running this refrigerator uses up energy, cutting into the fuel efficiency of the airplane, which is why typically only 50% of the cabin's air is replaced in each pumping cycle of 2-3 minutes. Although the airlines prefer to emphasize that this is a much faster recirculation rate than in the ventilation systems of most buildings, people are packed more tightly in an airplane.

2.1.4 Power

Power, P , is defined as the rate of change of energy, dE/dt . Power thus has units of joules per second, which are usually abbreviated as watts, 1 W=1 J/s. Since energy is conserved, we would have $dE/dt = 0$ if E was the total energy of a closed system, and that's not very interesting. What's usually more interesting to discuss is either the power flowing in or out of an open system, or the rate at which energy is being transformed from one form into another. The following is an example of energy flowing into an open system.

Heating by a lightbulb

example 4

- ▷ The electric company bills you for energy in units of kilowatt-hours (kilowatts multiplied by hours) rather than in SI units of joules. How many joules is a kilowatt-hour?
- ▷ $1 \text{ kilowatt-hour} = (1 \text{ kW})(1 \text{ hour}) = (1000 \text{ J/s})(3600 \text{ s}) = 3.6 \text{ MJ}$.

Now here's an example of energy being transformed from one form into another.

Human wattage

example 5

- ▷ Food contains chemical energy (discussed in more detail in section 2.4), and for historical reasons, food energy is normally given in non-SI units of Calories. One Calorie with a capital "C" equals 1000 calories, and 1 calorie is defined as 4.18 J. A typical person consumes 2000 Calories of food in a day, and converts nearly all of that directly to body heat. Compare the person's heat production to the rate of energy consumption of a 100-watt lightbulb.
- ▷ Strictly speaking, we can't really compute the derivative dE/dt ,

since we don't know how the person's metabolism ebbs and flows over the course of a day. What we can really compute is $\Delta E/\Delta t$, which is the power averaged over a one-day period.

Converting to joules, we find $\Delta E = 8 \times 10^6$ J for the amount of energy transformed into heat within our bodies in one day. Converting the time interval likewise into SI units, $\Delta t = 9 \times 10^4$ s. Dividing, we find that our power is $90 \text{ J/s} = 90 \text{ W}$, about the same as a lightbulb.

2.1.5 Gravitational energy

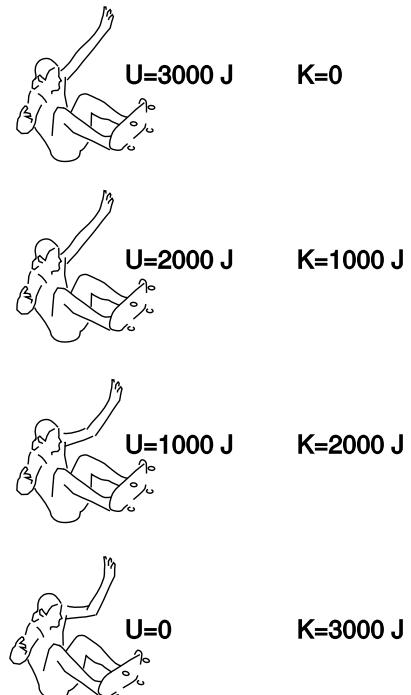
Gravitational energy, to which I've already alluded, is different from heat and kinetic energy in an important way. Heat and kinetic energy are properties of a single object, whereas gravitational energy describes an interaction between two objects. When the skater in figures g and h is at the top, his distance from the bulk of the planet earth is greater. Since we observe his kinetic energy decreasing on the way up, there must be some other form of energy that is increasing. We invent a new form of energy, called gravitational energy, and written U or U_g , which depends on the distance between his body and the planet. Where is this energy? It's not in the skater's body, and it's not inside the earth, either, since it takes two to tango. If either object didn't exist, there wouldn't be any interaction or any way to measure a distance, so it wouldn't make sense to talk about a distance-dependent energy. Just as marriage is a relationship between two people, gravitational energy is a relationship between two objects.

There is no precise way to define the distance between the skater and the earth, since both are objects that have finite size. As discussed in more detail in section 2.3, gravity is one of the fundamental forces of nature, a universal attraction between any two particles that have mass. Each atom in the skater's body is at a definite distance from each atom in the earth, but each of these distances is different. An atom in his foot is only a few centimeters from some of the atoms in the plaster side of the pool, but most of the earth's atoms are thousands of kilometers away from him. In theory, we might have to add up the contribution to the gravitational energy for every interaction between an atom in the skater's body and an atom in the earth.

For our present purposes, however, there is a far simpler and more practical way to solve problems. In any region of the earth's surface, there is a direction called "down," which we can establish by dropping a rock or hanging a plumb bob. In figure h, the skater is moving up and down in one dimension, and if we did measurements of his kinetic energy, like the made-up data in the figure, we could infer his gravitational energy. As long as we stay within a relatively small range of heights, we find that an object's gravitational energy increases at a steady rate with height. In other words, the strength



g / A skateboarder rises to the edge of an empty pool and then falls back down.



h / The sum of kinetic plus gravitational energy is constant.

of gravity doesn't change much if you only move up or down a few meters. We also find that the gravitational energy is proportional to the mass of the object we're testing. Writing y for the height, and g for the overall constant of proportionality, we have

$$U_g = mgy. \quad [\text{gravitational energy; } y=\text{height; only accurate within a small range of heights}]$$

The number g , with units of joules per kilogram per meter, is called the *gravitational field*. It tells us the strength of gravity in a certain region of space. Near the surface of our planet, it has a value of about $9.8 \text{ J/kg}\cdot\text{m}$, which is conveniently close to $10 \text{ J/kg}\cdot\text{m}$ for rough calculations.

Velocity at the bottom of a drop **example 6**

- ▷ If the skater in figure g drops 3 meters from rest, what is his velocity at the bottom of the pool?
- ▷ Starting from conservation of energy, we have

$$\begin{aligned} 0 &= \Delta E \\ &= \Delta K + \Delta U \\ &= K_f - K_i + U_f - U_i \\ &= \frac{1}{2}mv_f^2 + mgy_f - mgy_i && (\text{because } K_i=0) \\ &= \frac{1}{2}mv_f^2 + mg\Delta y, && (\Delta y < 0) \end{aligned}$$

so

$$\begin{aligned} v &= \sqrt{-2g\Delta y} \\ &= \sqrt{-(2)(10 \text{ J/kg}\cdot\text{m})(-3 \text{ m})} \\ &= 8 \text{ m/s} && (\text{rounded to one sig. fig.}) \end{aligned}$$

There are a couple of important things to note about this example. First, we were able to massage the equation so that it only involved Δy , rather than y itself. In other words, we don't need to worry about where $y = 0$ is; any coordinate system will work, as long as the positive y axis points up, not down. This is no accident. Gravitational energy can always be changed by adding a constant onto it, with no effect on the final result, as long as you're consistent within a given problem.

The other interesting thing is that the mass canceled out: even if the skater gained weight or strapped lead weights to himself, his velocity at the bottom would still be 8 m/s. This isn't an accident either. This is the same conclusion we reached in section 1.2, based on the equivalence of gravitational and inertial mass. The kinetic energy depends on the inertial mass, while gravitational energy is

related to gravitational mass, but since these two masses are equal, we were able to use a single symbol, m , for them, and cancel them out.

We can see from the equation $v = \sqrt{-2g\Delta y}$ that a falling object's velocity isn't constant. It increases as the object drops farther and farther. What about its acceleration? If we assume that air friction is negligible, the arguments in section 1.2 show that the acceleration can't depend on the object's mass, so there isn't much else the acceleration *can* depend on besides g . In fact, the acceleration of a falling object equals $-g$ (in a coordinate system where the positive y axis points up), as we can easily show using the chain rule:

$$\begin{aligned}\left(\frac{dv}{dt}\right) &= \left(\frac{dv}{dK}\right) \left(\frac{dK}{dU}\right) \left(\frac{dU}{dy}\right) \left(\frac{dy}{dt}\right) \\ &= \left(\frac{1}{mv}\right) (-1)(mg)(v) \\ &= -g,\end{aligned}$$

where I've calculated dv/dK as $1/(dK/dv)$, and $dK/dU = -1$ can be found by differentiating $K+U = (\text{constant})$ to give $dK+dU = 0$.⁷

We can also check that the units of g , $\text{J/kg}\cdot\text{m}$, are equivalent to the units of acceleration,

$$\frac{\text{J}}{\text{kg}\cdot\text{m}} = \frac{\text{kg}\cdot\text{m}^2/\text{s}^2}{\text{kg}\cdot\text{m}} = \frac{\text{m}}{\text{s}^2},$$

and therefore the strength of the gravitational field near the earth's surface can just as well be stated as 10 m/s^2 .

Speed after a given time

example 7

▷ An object falls from rest. How fast is it moving after two seconds? Assume that the amount of energy converted to heat by air friction is negligible.

▷ Under the stated assumption, we have $a = -g$, which can be integrated to give $v = -gt + \text{constant}$. If we let $t = 0$ be the beginning of the fall, then the constant of integration is zero, so at $t = 2 \text{ s}$ we have $v = -gt = -(10 \text{ m/s}^2) \times (2 \text{ s}) = 20 \text{ m/s}$.

The Vomit Comet

example 8

▷ The U.S. Air Force has an airplane, affectionately known as the Vomit Comet, in which astronaut trainees can experience simulated weightlessness. Oversimplifying a little, imagine that the plane climbs up high, and then drops straight down like a rock. (It actually flies along a parabola.) Since the people are falling with the same acceleration as the plane, the sensation is just like what you'd experience if you went out of the earth's gravitational

⁷There is a mathematical loophole in this argument that would allow the object to hover for a while with zero velocity and zero acceleration. This point is discussed on page 1025.

field. If the plane can start from 10 km up, what is the maximum amount of time for which the dive can last?

▷ Based on data about acceleration and distance, we want to find time. Acceleration is the second derivative of distance, so if we integrate the acceleration twice with respect to time, we can find how position relates to time. For convenience, let's pick a coordinate system in which the positive y axis is down, so $a=g$ instead of $-g$.

$$a = g$$

$$v = gt + \text{constant} \quad (\text{integrating})$$

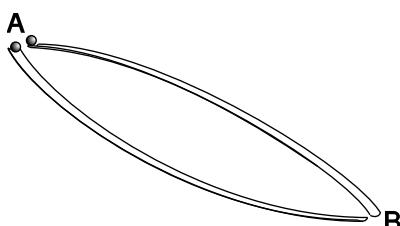
$$= gt \quad (\text{starts from rest})$$

$$y = \frac{1}{2}gt^2 + \text{constant} \quad (\text{integrating again})$$

Choosing our coordinate system to have $y = 0$ at $t = 0$, we can make the second constant of integration equal zero as well, so

$$\begin{aligned} t &= \sqrt{\frac{2y}{g}} \\ &= \sqrt{\frac{2 \cdot 10000 \text{ m}}{10 \text{ m/s}^2}} \\ &= \sqrt{2000 \text{ s}^2} \\ &= 40 \text{ s} \quad (\text{to one sig. fig.}) \end{aligned}$$

Note that if we hadn't converted the altitude to units of meters, we would have gotten the wrong answer, but we would have been alerted to the problem because the units inside the square root wouldn't have come out to be s^2 . In general, it's a good idea to convert all your data into SI (meter-kilogram-second) units before you do anything with them.



i / Two balls start from rest, and roll from A to B by different paths.

High road, low road

example 9

▷ In figure i, what can you say based on conservation of energy about the speeds of the balls when they reach point B? What does conservation of energy tell you about which ball will get there first? Assume friction doesn't convert any mechanical energy to heat or sound energy.

▷ Since friction is assumed to be negligible, there are only two forms of energy involved: kinetic and gravitational. Since both balls start from rest, and both lose the same amount of gravitational energy, they must have the same kinetic energy at the end, and therefore they're rolling at the same speed when they reach B. (A subtle point is that the balls have kinetic energy both because they're moving through space and because they're spinning as they roll. These two types of energy must be in fixed

proportion to one another, so this has no effect on the conclusion.)

Conservation of energy does not, however, tell us anything obvious about which ball gets there first. This is a general problem with applying conservation laws: conservation laws don't refer directly to time, since they are statements that something stays the same at all moments in time. We expect on intuitive grounds that the ball that goes by the lower ramp gets to B first, since it builds up speed early on.

Buoyancy

example 10

▷ A cubical box with mass m and volume $V = b^3$ is submerged in a fluid of density ρ . How much energy is required to raise it through a height Δy ?

▷ As the box moves up, it invades a volume $V' = b^2\Delta y$ previously occupied by some of the fluid, and fluid flows into an equal volume that it has vacated on the bottom. Lowering this amount of fluid by a height b reduces the fluid's gravitational energy by $\rho V'gb = \rho g b^3 \Delta y$, so the net change in energy is

$$\begin{aligned}\Delta E &= mg\Delta y - \rho g b^3 \Delta y \\ &= (m - \rho V)g\Delta y.\end{aligned}$$

In other words, it's as if the mass of the box had been reduced by an amount equal to the fluid that otherwise would have occupied that volume. This is known as Archimedes' principle, and it is true even if the box is not a cube, although we'll defer the more general proof until page 207 in chapter 3. If the box is less dense than the fluid, then it will float.

A simple machine

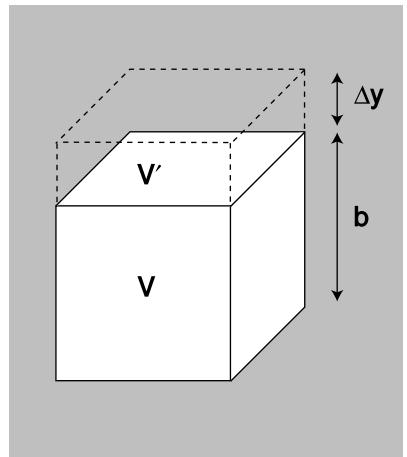
example 11

▷ If the father and son on the seesaw in figure k start from rest, what will happen?

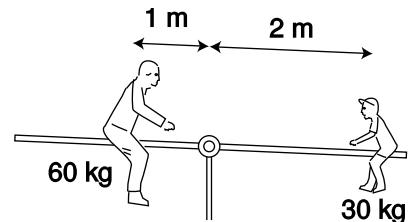
▷ Note that although the father is twice as massive, he is at half the distance from the fulcrum. If the seesaw was going to start rotating, it would have to be losing gravitational energy in order to gain some kinetic energy. However, there is no way for it to gain or lose gravitational energy by rotating in either direction. The change in gravitational energy would be

$$\begin{aligned}\Delta U &= \Delta U_1 + \Delta U_2 \\ &= g(m_1\Delta y_1 + m_2\Delta y_2),\end{aligned}$$

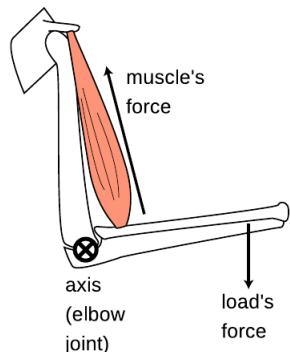
but Δy_1 and Δy_2 have opposite signs and are in the proportion of two to one, since the son moves along a circular arc that covers the same angle as the father's but has half the radius. Therefore $\Delta U = 0$, and there is no way for the seesaw to trade gravitational energy for kinetic.



j / How much energy is required to raise the submerged box through a height Δy ?



k / A seesaw.



l / The biceps muscle is a reversed lever.

The seesaw example demonstrates the principle of the lever, which is one of the basic mechanical building blocks known as simple machines. As discussed in more detail in chapters 3 and 4, the principle applies even when the interactions involved aren't gravitational.

Note that although a lever makes it easier to lift a heavy weight, it also decreases the distance traveled by the load. By reversing the lever, we can make the load travel a greater distance, at the expense of increasing the amount of force required. The human muscular-skeletal system uses reversed levers of this kind, which allows us to move more rapidly, and also makes our bodies more compact, at the expense of brute strength. A piano uses reversed levers so that a small amount of motion of the key produces a longer swing of the hammer. Another interesting example is the hydraulic jack shown in figure n. The analysis in terms of gravitational energy is exactly the same as for the seesaw, except that the relationship between Δy_1 and Δy_2 is now determined not by geometry but by conservation of mass: since water is highly incompressible, conservation of mass is approximately the same as a requirement of constant volume, which can only be satisfied if the distance traveled by each piston is in inverse proportion to its cross-sectional area.

Discussion Questions

A Hydroelectric power (water flowing over a dam to spin turbines) appears to be completely free. Does this violate conservation of energy? If not, then what is the ultimate source of the electrical energy produced by a hydroelectric plant?

B You throw a steel ball up in the air. How can you prove based on conservation of energy that it has the same speed when it falls back into your hand? What if you threw a feather up? Is energy not conserved in this case?

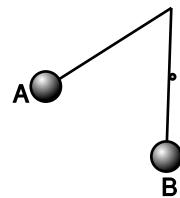
C Figure m shows a pendulum that is released at A and caught by a peg as it passes through the vertical, B. To what height will the bob rise on the right?

- D** What is wrong with the following definitions of g ?
- " g is gravity."
 - " g is the speed of a falling object."
 - " g is how hard gravity pulls on things."

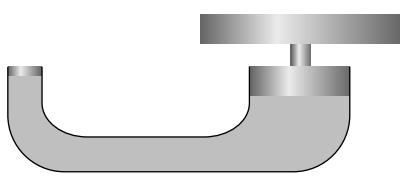
2.1.6 Equilibrium and stability

The seesaw in figure k is in equilibrium, meaning that if it starts out being at rest, it will stay put. This is known as a neutral equilibrium, since the seesaw has no preferred position to which it will return if we disturb it. If we move it to a different position and release it, it will stay at rest there as well. If we put it in motion, it will simply continue in motion until one person's feet hit the ground.

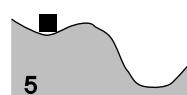
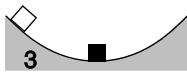
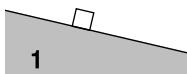
Most objects around you are in stable equilibria, like the black



m / Discussion question C.



n / A hydraulic jack.



o / The surfaces are frictionless. The black blocks are in equilibrium.

block in figure o/3. Even if the block is moved or set in motion, it will oscillate about the equilibrium position. The pictures are like graphs of y versus x , but since the gravitational energy $U = mgy$ is proportional to y , we can just as well think of them as graphs of U versus x . The block's stable equilibrium position is where the function $U(x)$ has a local minimum. The book you're reading right now is in equilibrium, but gravitational energy isn't the only form of energy involved. To move it upward, we'd have to supply gravitational energy, but downward motion would require a different kind of energy, in order to compress the table more. (As we'll see in section 2.4, this is electrical energy due to interactions between atoms within the table.)

A differentiable function's local extrema occur where its derivative is zero. A position where dU/dx is zero can be a stable (3), neutral (2), or unstable equilibrium, (4). An unstable equilibrium is like a pencil balanced on its tip. Although it could theoretically remain balanced there forever, in reality it will topple due to any tiny perturbation, such as an air current or a vibration from a passing truck. This is a technical, mathematical definition of instability, which is more restrictive than the way the word is used in ordinary speech. Most people would describe a domino standing upright as being unstable, but in technical usage it would be considered stable, because a certain finite amount of energy is required to tip it over, and perturbations smaller than that would only cause it to oscillate around its equilibrium position.

The domino is also an interesting example because it has two local minima, one in which it is upright, and another in which it is lying flat. A local minimum that is not the global minimum, as in figure o/5, is referred to as a metastable equilibrium.



p / Example 12.

A neutral equilibrium

example 12

Figure p shows a special-purpose one-block funicular railroad near Hill and Fourth Streets in Los Angeles, California, used for getting passengers up and down a very steep hill. It has two cars attached to a single loop of cable, arranged so that while one car goes up, the other comes down. They pass each other in the middle. Since one car's gravitational energy is increasing while

the other's is decreasing, the system is in neutral equilibrium. If there were no frictional heating, exactly zero energy would be required in order to operate the system. A similar counterweighting principle is used in aerial tramways in mountain resorts, and in elevators (with a solid weight, rather than a second car, as counterweight).

Water in a U-shaped tube

example 13

▷ The U-shaped tube in figure q has cross-sectional area A , and the density of the water inside is ρ . Find the gravitational energy as a function of the quantity y shown in the figure, and show that there is an equilibrium at $y=0$.

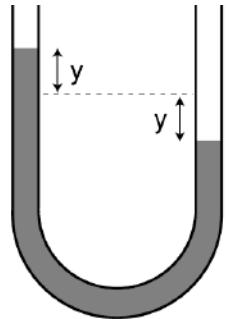
▷ The question is a little ambiguous, since gravitational energy is only well defined up to an additive constant. To fix this constant, let's define U to be zero when $y=0$. The difference between $U(y)$ and $U(0)$ is the energy that would be required to lift a water column of height y out of the right side, and place it above the dashed line, on the left side, raising it through a height y . This water column has height y and cross-sectional area A , so its volume is Ay , its mass is ρAy , and the energy required is $mgy = (\rho Ay)gy = \rho gAy^2$. We then have $U(y) = U(0) + \rho gAy^2 = \rho gAy^2$.

To find equilibria, we look for places where the derivative $dU/dy = 2\rho gAy$ equals 0. As we'd expect intuitively, the only equilibrium occurs at $y=0$. The second derivative test shows that this is a local minimum (not a maximum or a point of inflection), so this is a stable equilibrium.

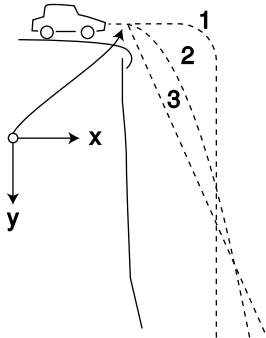
2.1.7 Predicting the direction of motion

Kinetic energy doesn't depend on the direction of motion. Sometimes this is helpful, as in the high road-low road example (p. 84, example 9), where we were able to predict that the balls would have the same final speeds, even though they followed different paths and were moving in different directions at the end. In general, however, the two conservation laws we've encountered so far aren't enough to predict an object's path through space, for which we need conservation of momentum (chapter 3), and the mathematical technique of vectors. Before we develop those ideas in their full generality, however, it will be helpful to do a couple of simple examples, including one that we'll get a lot of mileage out of in section 2.3.

Suppose we observe an air hockey puck gliding frictionlessly to the right at a velocity v , and we want to predict its future motion. Since there is no friction, no kinetic energy is converted to heat. The only form of energy involved is kinetic energy, so conservation of energy, $\Delta E = 0$, becomes simply $\Delta K = 0$. There's no particular reason for the puck to do anything but continue moving to the right at constant speed, but it would be equally consistent with conservation of energy if it spontaneously decided to reverse its direction of motion, changing its velocity to $-v$. Either way, we'd have $\Delta K = 0$. There is, however, a way to tell which motion is physical and which is unphysical. Suppose we consider the whole thing again in the frame of reference that is initially moving right along with the puck. In this frame, the puck starts out with $K = 0$. What we originally described as a reversal of its velocity from v to $-v$ is, in this



q / Water in a U-shaped tube.



r/A car drives over a cliff.

new frame of reference, a change from zero velocity to $-2v$, which would violate conservation of energy. In other words, the physically possible motion conserves energy in all frames of reference, but the unphysical motion only conserves energy in one special frame of reference.

For our second example, we consider a car driving off the edge of a cliff (r). For simplicity, we assume that air friction is negligible, so only kinetic and gravitational energy are involved. Does the car follow trajectory 1, familiar from Road Runner cartoons, trajectory 2, a parabola, or 3, a diagonal line? All three could be consistent with conservation of energy, in the ground's frame of reference. For instance, the car would have constant gravitational energy along the initial horizontal segment of trajectory 1, so during that time it would have to maintain constant kinetic energy as well. Only a parabola, however, is consistent with conservation of energy combined with Galilean relativity. Consider the frame of reference that is moving horizontally at the same speed as that with which the car went over the edge. In this frame of reference, the cliff slides out from under the initially motionless car. The car can't just hover for a while, so trajectory 1 is out. Repeating the same math as in example 8 on p. 83, we have

$$x^* = 0, \quad y^* = \frac{1}{2}gt^2$$

in this frame of reference, where the stars indicate coordinates measured in the moving frame of reference. These coordinates are related to the ground-fixed coordinates (x, y) by the equations

$$x = x^* + vt \quad \text{and} \quad y = y^*,$$

where v is the velocity of one frame with respect to the other. We therefore have

$$x = vt, \quad y = \frac{1}{2}gt^2,$$

in our original frame of reference. Eliminating t , we can see that this has the form of a parabola:

$$y = \frac{g}{2v^2}x^2.$$

self-check A

What would the car's motion be like in the * frame of reference if it followed trajectory 3?

▷ Answer, p. 1058

2.2 Numerical techniques

Engineering majors are a majority of the students in the kind of physics course for which this book is designed, so most likely you fall into that category. Although you surely recognize that physics is an important part of your training, if you've had any exposure to how engineers really work, you're probably skeptical about the flavor of problem-solving taught in most science courses. You realize that not very many practical engineering calculations fall into the narrow range of problems for which an exact solution can be calculated with a piece of paper and a sharp pencil. Real-life problems are usually complicated, and typically they need to be solved by number-crunching on a computer, although we can often gain insight by working simple approximations that have algebraic solutions. Not only is numerical problem-solving more useful in real life, it's also educational; as a beginning physics student, I really only felt like I understood projectile motion after I had worked it both ways, using algebra and then a computer program. (This was back in the days when 64 kilobytes of memory was considered a lot.)

In this section, we'll start by seeing how to apply numerical techniques to some simple problems for which we know the answer in "closed form," i.e., a single algebraic expression without any calculus or infinite sums. After that, we'll solve a problem that would have made you world-famous if you could have done it in the seventeenth century using paper and a quill pen! Before you continue, you should read Appendix 1 on page 1023 that introduces you to the Python programming language.

First let's solve the trivial problem of finding how much time it takes an object moving at speed v to travel a straight-line distance $dist$. This closed-form answer is, of course, $dist/v$, but the point is to introduce the techniques we can use to solve other problems of this type. The basic idea is to divide the distance up into n equal parts, and add up the times required to traverse all the parts. The following Python function does the job. Note that you shouldn't type in the line numbers on the left, and you don't need to type in the comments, either. I've omitted the prompts $>>>$ and \dots in order to save space.

```
1 import math
2 def time1(dist,v,n):
3     x=0                      # Initialize the position.
4     dx = dist/n                # Divide dist into n equal parts.
5     t=0                        # Initialize the time.
6     for i in range(n):
7         x = x+dx              # Change x.
8         dt=dx/v                # time=distance/speed
9         t=t+dt                  # Keep track of elapsed time.
10    return t
```

How long does it take to move 1 meter at a constant speed of 1 m/s?
If we do this,

```
>>> print(time1(1.0,1.0,10))      # dist, v, n
0.9999999999999989
```

Python produces the expected answer by dividing the distance into ten equal 0.1-meter segments, and adding up the ten 0.1-second times required to traverse each one. Since the object moves at constant speed, it doesn't even matter whether we set `n` to 10, 1, or a million:

```
>>> print(time1(1.0,1.0,1))      # dist, v, n
1.0
```

Now let's do an example where the answer isn't obvious to people who don't know calculus: how long does it take an object to fall through a height `h`, starting from rest? We know from example 8 on page 83 that the exact answer, found using calculus, is $\sqrt{2h/g}$. Let's see if we can reproduce that answer numerically. The main difference between this program and the previous one is that now the velocity isn't constant, so we need to update it as we go along. Conservation of energy gives $mgh = (1/2)mv^2 + mgy$ for the velocity v at height y , so $v = -\sqrt{2g(h-y)}$. (We choose the negative root because the object is moving down, and our coordinate system has the positive y axis pointing up.)

```
1  import math
2  def time2(h,n):
3      g=9.8                  # gravitational field
4      y=h                    # Initialize the height.
5      v=0                    # Initialize the velocity.
6      dy = -h/n               # Divide h into n equal parts.
7      t=0                    # Initialize the time.
8      for i in range(n):
9          y = y+dy            # Change y. (Note dy<0.)
10         v = -math.sqrt(2*g*(h-y))    # from cons. of energy
11         dt=dy/v             # dy and v are <0, so dt is >0
12         t=t+dt              # Keep track of elapsed time.
13     return t
14
```

For $h=1.0$ m, the closed-form result is $\sqrt{2 \cdot 1.0 \text{ m} / 9.8 \text{ m/s}^2} = 0.45$ s. With the drop split up into only 10 equal height intervals, the numerical technique provides a pretty lousy approximation:

```
>>> print(time2(1.0,10))      # h, n
0.35864270709233342
```

But by increasing `n` to ten thousand, we get an answer that's as close as we need, given the limited accuracy of the raw data:

```
>>> print(time2(1.0,10000))      # h, n
0.44846664060793945
```

A subtle point is that we changed `y` in line 9, and *then* on line 10 we calculated `v`, which depends on `y`. Since `y` is only changing by a ten-thousandth of a meter with each step, you might think this wouldn't make much of a difference, and you'd be almost right, except for one small problem: if we swapped lines 9 and 10, then the very first time through the loop, we'd have `v=0`, which would produce a division-by-zero error when we calculated `dt`! Actually what would make the most sense would be to calculate the velocity at height `y` and the velocity at height `y+dy` (recalling that `dy` is negative), average them together, and use that value of `y` to calculate the best estimate of the velocity between those two points. Since the acceleration is constant in the present example, this modification results in a program that gives an exact result even for `n=1`:

```
1  import math
2  def time3(h,n):
3      g=9.8
4      y=h
5      v=0
6      dy = -h/n
7      t=0
8      for i in range(n):
9          y_old = y
10         y = y+dy
11         v_old = math.sqrt(2*g*(h-y_old))
12         v = math.sqrt(2*g*(h-y))
13         v_avg = -(v_old+v)/2.
14         dt=dy/v_avg
15         t=t+dt
16     return t
17

>>> print(time3(1.0,1))      # h, n
0.45175395145262565
```

Now we're ready to attack a problem that challenged the best minds of Europe back in the days when there were no computers. In 1696, the mathematician Johann Bernoulli posed the following

famous question. Starting from rest, an object slides frictionlessly over a curve joining the point (a, b) to the point $(0, 0)$. Of all the possible shapes that such a curve could have, which one gets the object to its destination in the least possible time, and how much time does it take? The optimal curve is called the *brachistochrone*, from the Greek “short time.” The solution to the brachistochrone problem evaded Bernoulli himself, as well as Leibniz, who had been one of the inventors of calculus. The English physicist Isaac Newton, however, stayed up late one night after a day’s work running the royal mint, and, according to legend, produced an algebraic solution at four in the morning. He then published it anonymously, but Bernoulli is said to have remarked that when he read it, he knew instantly from the style that it was Newton — he could “tell the lion from the mark of his claw.”

Rather than attempting an exact algebraic solution, as Newton did, we’ll produce a numerical result for the shape of the curve and the minimum time, in the special case of $a=1.0$ m and $b=1.0$ m. Intuitively, we want to start with a fairly steep drop, since any speed we can build up at the start will help us throughout the rest of the motion. On the other hand, it’s possible to go too far with this idea: if we drop straight down for the whole vertical distance, and then do a right-angle turn to cover the horizontal distance, the resulting time of 0.68 s is quite a bit longer than the optimal result, the reason being that the path is unnecessarily long. There are infinitely many possible curves for which we could calculate the time, but let’s look at third-order polynomials,

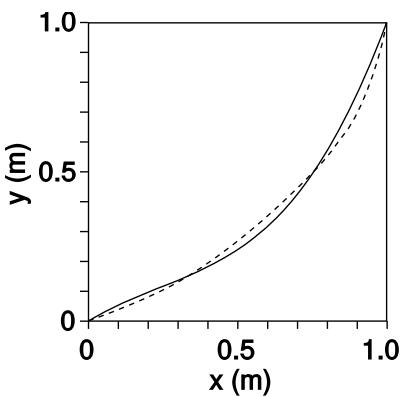
$$y = c_1x + c_2x^2 + c_3x^3,$$

where we require $c_3 = (b - c_1a - c_2a^2)/a^3$ in order to make the curve pass through the point (a, b) . The Python program, below, is not much different from what we’ve done before. The function only asks for c_1 and c_2 , and calculates c_3 internally at line 4. Since the motion is two-dimensional, we have to calculate the distance between one point and the next using the Pythagorean theorem, at line 16.

```

1 import math
2 def timeb(a,b,c1,c2,n):
3     g=9.8
4     c3 = (b-c1*a-c2*a**2)/(a**3)
5     x=a
6     y=b
7     dx = -a/n
8     t=0
9     for i in range(n):
10        y_old = y
11        x = x+dx
12        y = c1*x+c2*x**2+c3*x**3

```



a / Approximations to the brachistochrone curve using a third-order polynomial (solid line), and a seventh-order polynomial (dashed). The latter only improves the time by four milliseconds.

```

13     dy = y-y_old
14     v_old = math.sqrt(2*g*(b-y_old))
15     v = math.sqrt(2*g*(b-y))
16     v_avg = (v_old+v)/2.
17     ds = math.sqrt(dx**2+dy**2)           # Pythagorean thm.
18     dt=ds/v_avg
19     t=t+dt
20     return t
21

```

As a first guess, we could try a straight diagonal line, $y = x$, which corresponds to setting $c_1 = 1$, and all the other coefficients to zero. The result is a fairly long time:

```

>>> a=1.
>>> b=1.
>>> n=10000
>>> c1=1.
>>> c2=0.
>>> print(timeb(a,b,c1,c2,n))
0.63887656499994161

```

What we really need is a curve that's very steep on the right, and flatter on the left, so it would actually make more sense to try $y = x^3$:

```

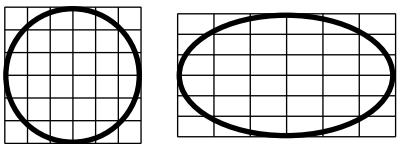
>>> c1=0.
>>> c2=0.
>>> print(timeb(a,b,c1,c2,n))
0.59458339947087069

```

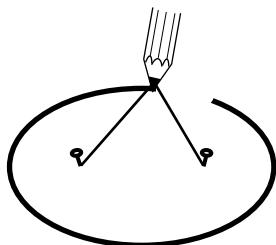
This is a significant improvement, and turns out to be only a hundredth of a second off of the shortest possible time! It's possible, although not very educational or entertaining, to find better approximations to the brachistochrone curve by fiddling around with the coefficients of the polynomial by hand. The real point of this discussion was to give an example of a nontrivial problem that can be attacked successfully with numerical techniques. I found the first approximation shown in figure a,

$$y = (0.62)x + (-0.93)x^2 + (1.31)x^3$$

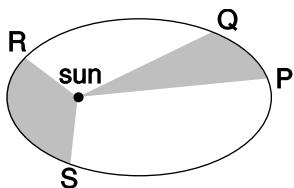
by using the program listed in appendix 2 on page 1026 to search automatically for the optimal curve. The seventh-order approximation shown in the figure came from a straightforward extension of the same program.



a / An ellipse is a circle that has been distorted by shrinking and stretching along perpendicular axes.



b / An ellipse can be constructed by tying a string to two pins and drawing like this with a pencil stretching the string taut. Each pin constitutes one focus of the ellipse.



c / If the time interval taken by the planet to move from P to Q is equal to the time interval from R to S, then according to Kepler's equal-area law, the two shaded areas are equal. The planet is moving faster during time interval RS than it was during PQ, because gravitational energy has been transformed into kinetic energy.

2.3 Gravitational phenomena

Cruise your radio dial today and try to find any popular song that would have been unimaginable without Louis Armstrong. By introducing solo improvisation into jazz, Armstrong took apart the jigsaw puzzle of popular music and fit the pieces back together in a different way. In the same way, Newton reassembled our view of the universe. Consider the titles of some recent physics books written for the general reader: **The God Particle, Dreams of a Final Theory**. When the subatomic particle called the neutrino was recently proven for the first time to have mass, specialists in cosmology began discussing seriously what effect this would have on calculations of the evolution of the universe from the Big Bang to its present state. Without the English physicist Isaac Newton, such attempts at universal understanding would not merely have seemed ambitious, they simply would not have occurred to anyone.

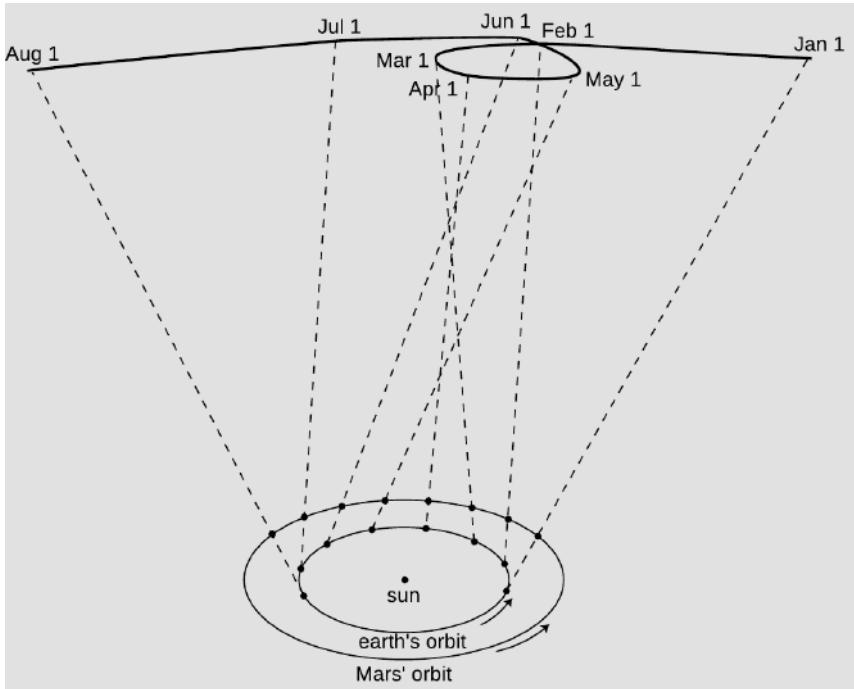
This section is about Newton's theory of gravity, which he used to explain the motion of the planets as they orbited the sun. Newton tosses off a general treatment of motion in the first 20 pages of his **Mathematical Principles of Natural Philosophy**, and then spends the next 130 discussing the motion of the planets. Clearly he saw this as the crucial scientific focus of his work. Why? Because in it he showed that the same laws of nature applied to the heavens as to the earth, and that the gravitational interaction that made an apple fall was the same as the one that kept the earth's motion from carrying it away from the sun.

2.3.1 Kepler's laws

Newton wouldn't have been able to figure out *why* the planets move the way they do if it hadn't been for the astronomer Tycho Brahe (1546-1601) and his protege Johannes Kepler (1571-1630), who together came up with the first simple and accurate description of *how* the planets actually do move. The difficulty of their task is suggested by the figure below, which shows how the relatively simple orbital motions of the earth and Mars combine so that as seen from earth Mars appears to be staggering in loops like a drunken sailor.

Brahe, the last of the great naked-eye astronomers, collected extensive data on the motions of the planets over a period of many years, taking the giant step from the previous observations' accuracy of about 10 minutes of arc (10/60 of a degree) to an unprecedented 1 minute. The quality of his work is all the more remarkable considering that his observatory consisted of four giant brass protractors mounted upright in his castle in Denmark. Four different observers would simultaneously measure the position of a planet in order to check for mistakes and reduce random errors.

With Brahe's death, it fell to his former assistant Kepler to try



d / As the earth and Mars revolve around the sun at different rates, the combined effect of their motions makes Mars appear to trace a strange, looped path across the background of the distant stars.

to make some sense out of the volumes of data. After 900 pages of calculations and many false starts and dead-end ideas, Kepler finally synthesized the data into the following three laws:

Kepler's elliptical orbit law: The planets orbit the sun in elliptical orbits with the sun at one focus.

Kepler's equal-area law: The line connecting a planet to the sun sweeps out equal areas in equal amounts of time.

Kepler's law of periods: The time required for a planet to orbit the sun, called its period, T , is proportional to the long axis of the ellipse raised to the $3/2$ power. The constant of proportionality is the same for all the planets.

Although the planets' orbits are ellipses rather than circles, most are very close to being circular. The earth's orbit, for instance, is only flattened by 1.7% relative to a circle. In the special case of a planet in a circular orbit, the two foci (plural of "focus") coincide at the center of the circle, and Kepler's elliptical orbit law thus says that the circle is centered on the sun. The equal-area law implies that a planet in a circular orbit moves around the sun with constant speed. For a circular orbit, the law of periods then amounts to a statement that the time for one orbit is proportional to $r^{3/2}$, where r is the radius. If all the planets were moving in their orbits at the same speed, then the time for one orbit would simply depend on the circumference of the circle, so it would only be proportional to r to the first power. The more drastic dependence on $r^{3/2}$ means that the outer planets must be moving more slowly than the inner planets.

Our main focus in this section will be to use the law of periods to deduce the general equation for gravitational energy. The equal-area law turns out to be a statement on conservation of angular momentum, which is discussed in chapter 4. We'll demonstrate the elliptical orbit law numerically in chapter 3, and analytically in chapter 4.

2.3.2 Circular orbits

Kepler's laws say that planets move along elliptical paths (with circles as a special case), which would seem to contradict the proof on page 90 that objects moving under the influence of gravity have parabolic trajectories. Kepler was right. The parabolic path was really only an approximation, based on the assumption that the gravitational field is constant, and that vertical lines are all parallel. In figure e, trajectory 1 is an ellipse, but it gets chopped off when the cannonball hits the earth, and the small piece of it that is above ground is nearly indistinguishable from a parabola. Our goal is to connect the previous calculation of parabolic trajectories, $y = (g/2v^2)x^2$, with Kepler's data for planets orbiting the sun in nearly circular orbits. Let's start by thinking in terms of an orbit that circles the earth, like orbit 2 in figure e. It's more natural now to choose a coordinate system with its origin at the center of the earth, so the parabolic approximation becomes $y = r - (g/2v^2)x^2$, where r is the distance from the center of the earth. For small values of x , i.e., when the cannonball hasn't traveled very far from the muzzle of the gun, the parabola is still a good approximation to the actual circular orbit, defined by the Pythagorean theorem, $r^2 = x^2 + y^2$, or $y = r\sqrt{1 - x^2/r^2}$. For small values of x , we can use the approximation $\sqrt{1 + \epsilon} \approx 1 + \epsilon/2$ to find $y \approx r - (1/2r)x^2$. Setting this equal to the equation of the parabola, we have $g/2v^2 = (1/2r)$, or

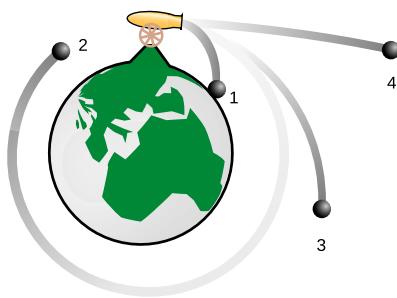
$$v = \sqrt{gr} \quad [\text{condition for a circular orbit}].$$

Low-earth orbit

example 14

To get a feel for what this all means, let's calculate the velocity required for a satellite in a circular low-earth orbit. Real low-earth-orbit satellites are only a few hundred km up, so for purposes of rough estimation we can take r to be the radius of the earth, and g is not much less than its value on the earth's surface, 10 m/s^2 . Taking numerical data from Appendix 4, we have

$$\begin{aligned} v &= \sqrt{gr} \\ &= \sqrt{(10 \text{ m/s}^2)(6.4 \times 10^3 \text{ km})} \\ &= \sqrt{(10 \text{ m/s}^2)(6.4 \times 10^6 \text{ m})} \\ &= \sqrt{6.4 \times 10^7 \text{ m}^2/\text{s}^2} \\ &= 8000 \text{ m/s} \end{aligned}$$



e / A cannon fires cannonballs at different velocities, from the top of an imaginary mountain that rises above the earth's atmosphere. This is almost the same as a figure Newton included in his **Mathematical Principles**.

(about twenty times the speed of sound).

In one second, the satellite moves 8000 m horizontally. During this time, it drops the same distance any other object would: about 5 m. But a drop of 5 m over a horizontal distance of 8000 m is just enough to keep it at the same altitude above the earth's curved surface.

2.3.3 The sun's gravitational field

We can now use the circular orbit condition $v = \sqrt{gr}$, combined with Kepler's law of periods, $T \propto r^{3/2}$ for circular orbits, to determine how the sun's gravitational field falls off with distance.⁸ From there, it will be just a hop, skip, and a jump to get to a universal description of gravitational interactions.

The velocity of a planet in a circular orbit is proportional to r/T , so

$$\begin{aligned}r/T &\propto \sqrt{gr} \\r/r^{3/2} &\propto \sqrt{gr} \\g &\propto 1/r^2\end{aligned}$$

If gravity behaves systematically, then we can expect the same to be true for the gravitational field created by any object, not just the sun.

There is a subtle point here, which is that so far, r has just meant the radius of a circular orbit, but what we have come up with smells more like an equation that tells us the strength of the gravitational field made by some object (the sun) if we know how far we are from the object. In other words, we could reinterpret r as the distance from the sun.

2.3.4 Gravitational energy in general

We now want to find an equation for the gravitational energy of any two masses that attract each other from a distance r . We assume that r is large enough compared to the distance between the objects so that we don't really have to worry about whether r is measured from center to center or in some other way. This would be a good approximation for describing the solar system, for example, since the sun and planets are small compared to the distances between them — that's why you see Venus (the “evening star”) with your bare eyes as a dot, not a disk.

The equation we seek is going to give the gravitational energy, U , as a function of m_1 , m_2 , and r . We already know from expe-

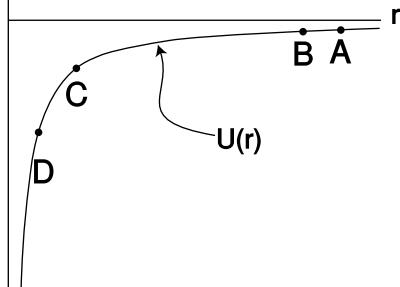
⁸There is a hidden assumption here, which is that the sun doesn't move. Actually the sun wobbles a little because of the planets' gravitational interactions with it, but the wobble is small due to the sun's large mass, so it's a pretty good approximation to assume the sun is stationary. Chapter 3 provides the tools to analyze this sort of thing completely correctly — see p. 144.

rience with gravity near the earth's surface that U is proportional to the mass of the object that interacts with the earth gravitationally, so it makes sense to assume the relationship is symmetric: U is presumably proportional to the product m_1m_2 . We can no longer assume $\Delta U \propto \Delta r$, as in the earth's-surface equation $\Delta U = mg\Delta y$, since we are trying to construct an equation that would be valid for all values of r , and g depends on r . We can, however, consider an infinitesimally small change in distance dr , for which we'll have $dU = m_2g_1 dr$, where g_1 is the gravitational field created by m_1 . (We could just as well have written this as $dU = m_1g_2 dr$, since we're not assuming either mass is "special" or "active.") Integrating this equation, we have

$$\begin{aligned} \int dU &= \int m_2g_1 dr \\ U &= m_2 \int g_1 dr \\ U &\propto m_1m_2 \int \frac{1}{r^2} dr \\ U &\propto -\frac{m_1m_2}{r}, \end{aligned}$$

where we're free to take the constant of integration to be equal to zero, since gravitational energy is never a well-defined quantity in absolute terms. Writing G for the constant of proportionality, we have the following fundamental description of gravitational interactions:

$$U = -\frac{Gm_1m_2}{r} \quad \begin{array}{l} \text{[gravitational energy of two masses} \\ \text{separated by a distance } r] \end{array}$$



f / The gravitational energy $U = -Gm_1m_2/r$ graphed as a function of r .

We'll refer to this as Newton's law of gravity, although in reality he stated it in an entirely different form, which turns out to be mathematically equivalent to this one.

Let's interpret his result. First, don't get hung up on the fact that it's negative, since it's only differences in gravitational energy that have physical significance. The graph in figure f could be shifted up or down without having any physical effect. The slope of this graph relates to the strength of the gravitational field. For instance, suppose figure f is a graph of the gravitational energy of an asteroid interacting with the sun. If the asteroid drops straight toward the sun, from A to B, the decrease in gravitational energy is very small, so it won't speed up very much during that motion. Points C and D, however, are in a region where the graph's slope is much greater. As the asteroid moves from C to D, it loses a lot of gravitational energy, and therefore speeds up considerably. This is due to the stronger gravitational field.

Determining G

example 15

The constant G is not easy to determine, and Newton went to his grave without knowing an accurate value for it. If we knew the mass of the earth, then we could easily determine G from experiments with terrestrial gravity, but the only way to determine the mass of the earth accurately in units of kilograms is by finding G and reasoning the other way around! (If you estimate the average density of the earth, you can make at least a rough estimate of G .) Figures g and h show how G was first measured by Henry Cavendish in the nineteenth century. The rotating arm is released from rest, and the kinetic energy of the two moving balls is measured when they pass position C. Conservation of energy gives

$$-2\frac{GMm}{r_{BA}} - 2\frac{GMm}{r_{BD}} = -2\frac{GMm}{r_{CA}} - 2\frac{GMm}{r_{CD}} + 2K,$$

where M is the mass of one of the large balls, m is the mass of one of the small ones, and the factors of two, which will cancel, occur because every energy is mirrored on the opposite side of the apparatus. (As discussed on page 102, it turns out that we get the right result by measuring all the distances from the center of one sphere to the center of the other.) This can easily be solved for G . The best modern value of G , from later versions of the same experiment, is $6.67 \times 10^{-11} \text{ J} \cdot \text{m/kg}^2$.

Escape velocity

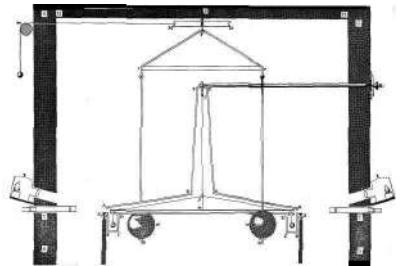
example 16

▷ The Pioneer 10 space probe was launched in 1972, and continued sending back signals for 30 years. In the year 2001, not long before contact with the probe was lost, it was about $1.2 \times 10^{13} \text{ m}$ from the sun, and was moving almost directly away from the sun at a velocity of $1.21 \times 10^4 \text{ m/s}$. The mass of the sun is $1.99 \times 10^{30} \text{ kg}$. Will Pioneer 10 escape permanently, or will it fall back into the solar system?

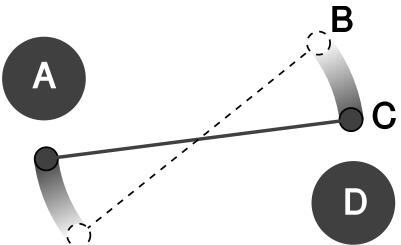
▷ We want to know whether there will be a point where the probe will turn around. If so, then it will have zero kinetic energy at the turnaround point:

$$\begin{aligned} K_i + U_i &= U_f \\ \frac{1}{2}mv^2 - \frac{GMm}{r_i} &= -\frac{GMm}{r_f} \\ \frac{1}{2}v^2 - \frac{GM}{r_i} &= -\frac{GM}{r_f}, \end{aligned}$$

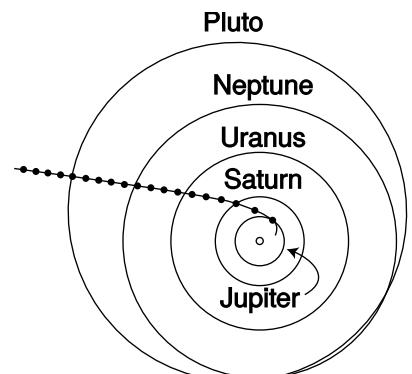
where M is the mass of the sun, m is the (irrelevant) mass of the probe, and r_f is the distance from the sun of the hypothetical turnaround point. Plugging in numbers on the left, we get a positive result. There can therefore be no solution, since the right side is negative. There won't be any turnaround point, and Pioneer 10 is never coming back.



g / Cavendish's original drawing of the apparatus for his experiment, discussed in example 15. The room was sealed to exclude air currents, and the motion was observed through telescopes sticking through holes in the walls.



h / A simplified drawing of the Cavendish experiment, viewed from above. The rod with the two small masses on the ends hangs from a thin fiber, and is free to rotate.



i / The Pioneer 10 space probe's trajectory from 1974 to 1992, with circles marking its position at one-year intervals. After its 1974 slingshot maneuver around Jupiter, the probe's motion was determined almost exclusively by the sun's gravity.

The minimum velocity required for this to happen is called *escape velocity*. For speeds above escape velocity, the orbits are open-ended hyperbolas, rather than repeating elliptical orbits. In figure i, Pioneer's hyperbolic trajectory becomes almost indistinguishable from a line at large distances from the sun. The motion slows perceptibly in the first few years after 1974, but later the speed becomes nearly constant, as shown by the nearly constant spacing of the dots.

The gravitational field

We got the energy equation $U = -Gm_1m_2/r$ by integrating $g \propto 1/r^2$ and then inserting a constant of proportionality to make the proportionality into an equation. The opposite of an integral is a derivative, so we can now go backwards and insert a constant of proportionality in $g \propto 1/r^2$ that will be consistent with the energy equation:

$$\begin{aligned} dU &= m_2 g_1 dr \\ g_1 &= \frac{1}{m_2} \frac{dU}{dr} \\ &= \frac{1}{m_2} \frac{d}{dr} \left(-\frac{Gm_1m_2}{r} \right) \\ &= -Gm_1 \frac{d}{dr} \left(\frac{1}{r} \right) \\ &= \frac{Gm_1}{r^2} \end{aligned}$$

This kind of inverse-square law occurs all the time in nature. For instance, if you go twice as far away from a lightbulb, you receive 1/4 as much light from it, because as the light spreads out, it is like an expanding sphere, and a sphere with twice the radius has four times the surface area. It's like spreading the same amount of peanut butter on four pieces of bread instead of one — we have to spread it thinner.

Discussion Questions

A A bowling ball interacts gravitationally with the earth. Would it make sense for the gravitational energy to be inversely proportional to the distance between their surfaces rather than their centers?

2.3.5 The shell theorem

Newton's great insight was that gravity near the earth's surface was the same kind of interaction as the one that kept the planets from flying away from the sun. He told his niece that the idea came to him when he saw an apple fall from a tree, which made him wonder whether the earth might be affecting the apple and the moon in the same way. Up until now, we've generally been dealing with gravitational interactions between objects that are small compared to the distances between them, but that assumption doesn't apply to

the apple. A kilogram of dirt a few feet under his garden in England would interact much more strongly with the apple than a kilogram of molten rock deep under Australia, thousands of miles away. Also, we know that the earth has some parts that are more dense, and some parts that are less dense. The solid crust, on which we live, is considerably less dense than the molten rock on which it floats. By all rights, the computation of the total gravitational energy of the apple should be a horrendous mess. Surprisingly, it turns out to be fairly simple in the end. First, we note that although the earth doesn't have the same density throughout, it does have spherical symmetry: if we imagine dividing it up into thin concentric shells, the density of each shell is uniform.

Second, it turns out that a uniform spherical shell interacts with external masses as if all its mass were concentrated at its center.

The shell theorem: The gravitational energy of a uniform spherical shell of mass M interacting with a pointlike mass m outside it equals $-GMm/s$, where s is the center-to-center distance. If mass m is inside the shell, then the energy is constant, i.e., the shell's interior gravitational field is zero.

Proof: Let b be the radius of the shell, h its thickness, and ρ its density. Its volume is then $V=(\text{area})(\text{thickness})=4\pi b^2 h$, and its mass is $M = \rho V = 4\pi \rho b^2 h$. The strategy is to divide the shell up into rings as shown in figure j, with each ring extending from θ to $\theta + d\theta$. Since the ring is infinitesimally skinny, its entire mass lies at the same distance, r , from mass m . The width of such a ring is found by the definition of radian measure to be $w = b d\theta$, and its mass is $dM = (\rho)(\text{circumference})(\text{thickness})(\text{width}) = (\rho)(2\pi b \sin \theta)(h)(b d\theta) = 2\pi \rho b^2 h \sin \theta d\theta$. The gravitational energy of the ring interacting with mass m is therefore

$$\begin{aligned} dU &= -\frac{Gm dM}{r} \\ &= -2\pi G\rho b^2 hm \frac{\sin \theta d\theta}{r}. \end{aligned}$$

Integrating both sides, we find the total gravitational energy of the shell:

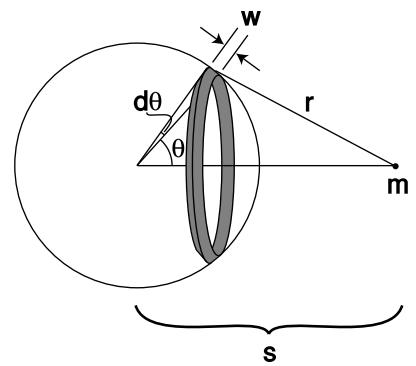
$$U = -2\pi G\rho b^2 hm \int_0^\pi \frac{\sin \theta d\theta}{r}$$

The integral has a mixture of the variables r and θ , which are related by the law of cosines,

$$r^2 = b^2 + s^2 - 2bs \cos \theta,$$

and to evaluate the integral, we need to get everything in terms of either r and dr or θ and $d\theta$. The relationship between the differentials is found by differentiating the law of cosines,

$$2r dr = 2bs \sin \theta d\theta,$$



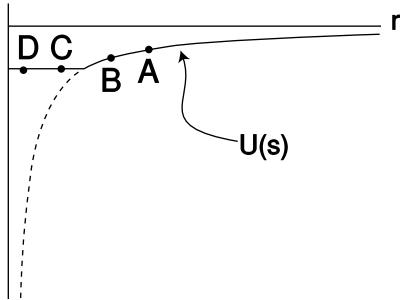
j / A spherical shell of mass M interacts with a pointlike mass m .

and since $\sin \theta d\theta$ occurs in the integral, the easiest path is to substitute for it, and get everything in terms of r and dr :

$$\begin{aligned} U &= -\frac{2\pi G\rho b hm}{s} \int_{s-b}^{s+b} dr \\ &= -\frac{4\pi G\rho b^2 hm}{s} \\ &= -\frac{GMm}{s} \end{aligned}$$

This was all under the assumption that mass m was on the outside of the shell. To complete the proof, we consider the case where it's inside. In this case, the only change is that the limits of integration are different:

$$\begin{aligned} U &= -\frac{2\pi G\rho b hm}{s} \int_{b-s}^{b+s} dr \\ &= -4\pi G\rho b hm \\ &= -\frac{GMm}{b} \end{aligned}$$



k / The gravitational energy of a mass m at a distance s from the center of a hollow spherical shell of mass M .

The two results are equal at the surface of the sphere, $s = b$, so the constant-energy part joins continuously onto the $1/s$ part, and the effect is to chop off the steepest part of the graph that we would have had if the whole mass M had been concentrated at its center. Dropping a mass m from A to B in figure k releases the same amount of energy as if mass M had been concentrated at its center, but there is no release of gravitational energy at all when moving between two interior points like C and D. In other words, the internal gravitational field is zero. Moving from C to D brings mass m farther away from the nearby side of the shell, but closer to the far side, and the cancellation between these two effects turns out to be perfect. Although the gravitational field has to be zero at the center due to symmetry, it's much more surprising that it cancels out perfectly in the whole interior region; this is a special mathematical characteristic of a $1/r$ interaction like gravity.

Newton's apple

example 17

Over a period of 27.3 days, the moon travels the circumference of its orbit, so using data from Appendix 4, we can calculate its speed, and solve the circular orbit condition to determine the strength of the earth's gravitational field at the moon's distance from the earth, $g = v^2/r = 2.72 \times 10^{-3} \text{ m/s}^2$, which is 3600 times smaller than the gravitational field at the earth's surface. The center-to-center distance from the moon to the earth is 60 times greater than the radius of the earth. The earth is, to a very good approximation, a sphere made up of concentric shells, each with uniform density, so the shell theorem tells us that its external gravitational field is the same as if all its mass was concentrated

at its center. We already know that a gravitational energy that varies as $-1/r$ is equivalent to a gravitational field proportional to $1/r^2$, so it makes sense that a distance that is greater by a factor of 60 corresponds to a gravitational field that is $60 \times 60 = 3600$ times weaker. Note that the calculation didn't require knowledge of the earth's mass or the gravitational constant, which Newton didn't know.

In 1665, shortly after Newton graduated from Cambridge, the Great Plague forced the college to close for two years, and Newton returned to the family farm and worked intensely on scientific problems. During this productive period, he carried out this calculation, but it came out wrong, causing him to doubt his new theory of gravity. The problem was that during the plague years, he was unable to use the university's library, so he had to use a figure for the radius of the moon's orbit that he had memorized, and he forgot that the memorized value was in units of nautical miles rather than statute miles. Once he realized his mistake, he found that the calculation came out just right, and became confident that his theory was right after all.⁹

Weighing the earth

example 18

▷ Once Cavendish had found $G = 6.67 \times 10^{-11} \text{ J} \cdot \text{m/kg}^2$ (p. 101, example 15), it became possible to determine the mass of the earth. By the shell theorem, the gravitational energy of a mass m at a distance r from the center of the earth is $U = -GMm/r$, where M is the mass of the earth. The gravitational field is related to this by $mg dr = dU$, or $g = (1/m) dU/dr = GM/r^2$. Solving for M , we have

$$\begin{aligned} M &= gr^2/G \\ &= \frac{(9.8 \text{ m/s}^2)(6.4 \times 10^6 \text{ m})^2}{6.67 \times 10^{-11} \text{ J} \cdot \text{m/kg}^2} \\ &= 6.0 \times 10^{24} \frac{\text{m}^2 \cdot \text{kg}^2}{\text{J} \cdot \text{s}^2} \\ &= 6.0 \times 10^{24} \text{ kg} \end{aligned}$$

Gravity inside the earth

example 19

▷ The earth is somewhat more dense at greater depths, but as an approximation let's assume it has a constant density throughout. How does its internal gravitational field vary with the distance r from the center?

▷ Let's write b for the radius of the earth. The shell theorem tell us that at a given location r , we only need to consider the mass $M_{<r}$

⁹Some historians are suspicious that the story of the apple and the mistake in conversions may have been fabricated by Newton later in life. The conversion incident may have been a way of explaining his long delay in publishing his work, which led to a conflict with Leibniz over priority in the invention of calculus.

that is deeper than r . Under the assumption of constant density, this mass is related to the total mass of the earth by

$$\frac{M_{\leq r}}{M} = \frac{r^3}{b^3},$$

and by the same reasoning as in example 18,

$$g = \frac{GM_{\leq r}}{r^2},$$

so

$$g = \frac{GMr}{b^3}.$$

In other words, the gravitational field interpolates linearly between zero at $r = 0$ and its ordinary surface value at $r = b$.

The following example applies the numerical techniques of section 2.2.

From the earth to the moon

example 20

The Apollo 11 mission landed the first humans on the moon in 1969. In this example, we'll estimate the time it took to get to the moon, and compare our estimate with the actual time, which was 73.0708 hours from the engine burn that took the ship out of earth orbit to the engine burn that inserted it into lunar orbit. During this time, the ship was coasting with the engines off, except for a small course-correction burn, which we neglect. More importantly, we do the calculation for a straight-line trajectory rather than the real S-shaped one, so the result can only be expected to agree roughly with what really happened. The following data come from the original press kit, which NASA has scanned and posted on the Web:

$$\begin{aligned} \text{initial altitude} &= 3.363 \times 10^5 \text{ m} \\ \text{initial velocity} &= 1.083 \times 10^4 \text{ m/s} \end{aligned}$$

The endpoint of the straight-line trajectory is a free-fall impact on the lunar surface, which is also unrealistic (luckily for the astronauts).

The ship's energy is

$$E = -\frac{GM_e m}{r} - \frac{GM_m m}{r_m - r} + \frac{1}{2}mv^2,$$

but since everything is proportional to the mass of the ship, m , we can divide it out

$$\frac{E}{m} = -\frac{GM_e}{r} - \frac{GM_m}{r_m - r} + \frac{1}{2}v^2,$$

and the energy variables in the program with names like e , k , and u are actually energies per unit mass. The program is a straightforward modification of the function `time3` on page 93.

```

1 import math
2 def tmoon(vi,ri,rf,n):
3     bigg=6.67e-11      # gravitational constant
4     me=5.97e24         # mass of earth
5     mm=7.35e22         # mass of moon
6     rm=3.84e8          # earth-moon distance
7     r=ri
8     v=vi
9     dr = (rf-ri)/n
10    e=-bigg*me/ri-bigg*mm/(rm-ri)+.5*vi**2
11    t=0
12    for i in range(n):
13        u_old = -bigg*me/r-bigg*mm/(rm-r)
14        k_old = e - u_old
15        v_old = math.sqrt(2.*k_old)
16        r = r+dr
17        u = -bigg*me/r-bigg*mm/(rm-r)
18        k = e - u
19        v = math.sqrt(2.*k)
20        v_avg = .5*(v_old+v)
21        dt=dr/v_avg
22        t=t+dt
23    return t
24

>>> re=6.378e6 # radius of earth
>>> rm=1.74e6 # radius of moon
>>> ri=re+3.363e5 # re+initial altitude
>>> rf=3.8e8-rm # earth-moon distance minus rm
>>> vi=1.083e4 # initial velocity
>>> print(tmoon(vi,ri,rf,1000)/3600.) # convert seconds to hours
59.654047441976552

```

This is pretty decent agreement, considering the wildly inaccurate trajectory assumed. It's interesting to see how much the duration of the trip changes if we increase the initial velocity by only ten percent:

```

>>> vi=1.2e4
>>> print(tmoon(vi,ri,rf,1000)/3600.)
18.177752636111677

```

The most important reason for using the lower speed was that if something had gone wrong, the ship would have been able to whip around the moon and take a "free return" trajectory back to the earth, without having to do any further burns. At a higher speed, the ship would have had so much kinetic energy that in the absence of any further engine burns, it would have escaped

from the earth-moon system. The Apollo 13 mission had to take a free return trajectory after an explosion crippled the spacecraft.

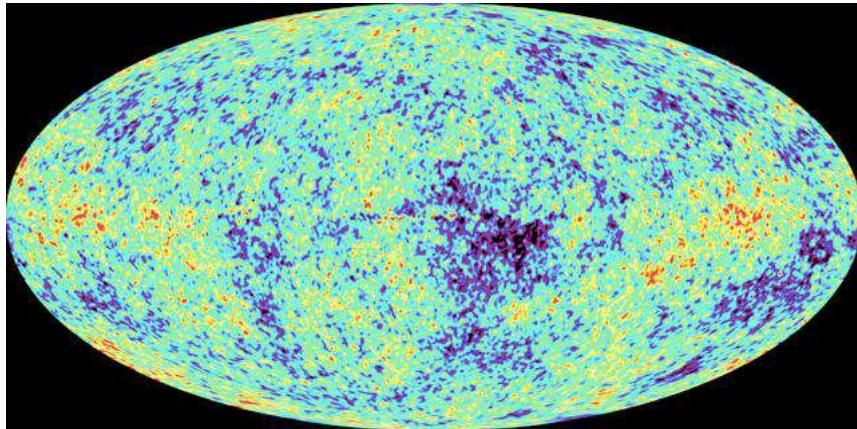
2.3.6 * Evidence for repulsive gravity

Until recently, physicists thought they understood gravity fairly well. Einstein had modified Newton's theory, but certain characteristics of gravitational forces were firmly established. For one thing, they were always attractive. If gravity always attracts, then it is logical to ask why the universe doesn't collapse. Newton had answered this question by saying that if the universe was infinite in all directions, then it would have no geometric center toward which it would collapse; the forces on any particular star or planet exerted by distant parts of the universe would tend to cancel out by symmetry. More careful calculations, however, show that Newton's universe would have a tendency to collapse on smaller scales: any part of the universe that happened to be slightly more dense than average would contract further, and this contraction would result in stronger gravitational forces, which would cause even more rapid contraction, and so on.

When Einstein overhauled gravity, the same problem reared its ugly head. Like Newton, Einstein was predisposed to believe in a universe that was static, so he added a special repulsive term to his equations, intended to prevent a collapse. This term was not associated with any interaction of mass with mass, but represented merely an overall tendency for space itself to expand unless restrained by the matter that inhabited it. It turns out that Einstein's solution, like Newton's, is unstable. Furthermore, it was soon discovered observationally that the universe was expanding, and this was interpreted by creating the Big Bang model, in which the universe's current expansion is the aftermath of a fantastically hot explosion.¹⁰ An expanding universe, unlike a static one, was capable of being explained with Einstein's equations, without any repulsion term. The universe's expansion would simply slow down over time due to the attractive gravitational forces. After these developments, Einstein said woefully that adding the repulsive term, known as the cosmological constant, had been the greatest blunder of his life.

This was the state of things until 1999, when evidence began to turn up that the universe's expansion has been speeding up rather than slowing down! The first evidence came from using a telescope as a sort of time machine: light from a distant galaxy may have taken billions of years to reach us, so we are seeing it as it was far in the past. Looking back in time, astronomers saw the universe expanding at speeds that were lower, rather than higher. At first they were mortified, since this was exactly the opposite of what had been expected. The statistical quality of the data was also not good enough to constitute ironclad proof, and there were worries

¹⁰Subsection 6.1.5 presents some evidence for the Big Bang theory.



m / The WMAP probe's map of the cosmic microwave background is like a "baby picture" of the universe.

about systematic errors. The case for an accelerating expansion has however been supported by high-precision mapping of the dim, sky-wide afterglow of the Big Bang, known as the cosmic microwave background.

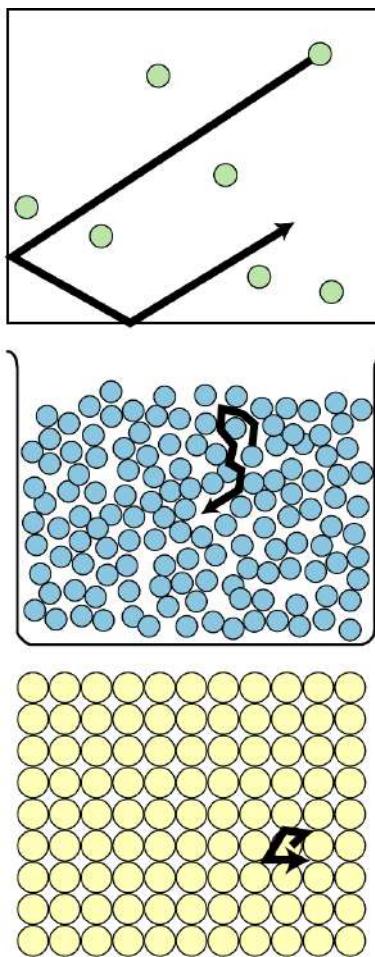
So now Einstein's "greatest blunder" has been resurrected. Since we don't actually know whether or not this self-repulsion of space has a constant strength, the term "cosmological *constant*" has lost currency. Nowadays physicists usually refer to the phenomenon as "dark energy." Picking an impressive-sounding name for it should not obscure the fact that we know absolutely nothing about the nature of the effect or why it exists.

2.4 Atomic phenomena

Variety is the spice of life, not of science. So far this chapter has focused on heat energy, kinetic energy, and gravitational energy, but it might seem that in addition to these there is a bewildering array of other forms of energy. Gasoline, chocolate bars, batteries, melting water — in each case there seems to be a whole new type of energy. The physicist's psyche rebels against the prospect of a long laundry list of types of energy, each of which would require its own equations, concepts, notation, and terminology. The point at which we've arrived in the study of energy is analogous to the period in the 1960's when a half a dozen new subatomic particles were being discovered every year in particle accelerators. It was an embarrassment. Physicists began to speak of the "particle zoo," and it seemed that the subatomic world was distressingly complex. The particle zoo was simplified by the realization that most of the new particles being whipped up were simply clusters of a previously unsuspected set of fundamental particles (which were whimsically dubbed quarks, a made-up word from a line of poetry by James Joyce, "Three quarks for Master Mark.") The energy zoo can also be simplified, and it's the purpose of this section to demonstrate the hidden similarities between forms of energy as seemingly different



a / A vivid demonstration that heat is a form of motion. A small amount of boiling water is poured into the empty can, which rapidly fills up with hot steam. The can is then sealed tightly, and soon crumples.



b / Random motion of atoms in a gas, a liquid, and a solid.

as heat and motion.

2.4.1 Heat is kinetic energy.

What is heat really? Is it an invisible fluid that your bare feet soak up from a hot sidewalk? Can one ever remove all the heat from an object? Is there a maximum to the temperature scale?

The theory of heat as a fluid seemed to explain why colder objects absorbed heat from hotter ones, but once it became clear that heat was a form of energy, it began to seem unlikely that a material substance could transform itself into and out of all those other forms of energy like motion or light. For instance, a compost pile gets hot, and we describe this as a case where, through the action of bacteria, chemical energy stored in the plant cuttings is transformed into heat energy. The heating occurs even if there is no nearby warmer object that could have been leaking “heat fluid” into the pile.

An alternative interpretation of heat was suggested by the theory that matter is made of atoms. Since gases are thousands of times less dense than solids or liquids, the atoms (or clusters of atoms called molecules) in a gas must be far apart. In that case, what is keeping all the air molecules from settling into a thin film on the floor of the room in which you are reading this book? The simplest explanation is that they are moving very rapidly, continually ricocheting off of the floor, walls, and ceiling. Though bizarre, the cloud-of-bullets image of a gas did give a natural explanation for the surprising ability of something as tenuous as a gas to exert huge forces.

The experiment shown in figure a, for instance, can be explained as follows. The high temperature of the steam is interpreted as a high average speed of random motions of its molecules. Before the lid was put on the can, the rapidly moving steam molecules pushed their way out of the can, forcing the slower air molecules out of the way. As the steam inside the can thinned out, a stable situation was soon achieved, in which the force from the less dense steam molecules moving at high speed balanced against the force from the more dense but slower air molecules outside. The cap was put on, and after a while the steam inside the can began to cool off. The force from the cooler, thin steam no longer matched the force from the cool, dense air outside, and the imbalance of forces crushed the can.

This type of observation leads naturally to the conclusion that hotter matter differs from colder in that its atoms' random motion is more rapid. In a liquid, the motion could be visualized as people in a milling crowd shoving past each other more quickly. In a solid, where the atoms are packed together, the motion is a random vibration of each atom as it knocks against its neighbors.

We thus achieve a great simplification in the theory of heat. Heat is simply a form of kinetic energy, the total kinetic energy of random

motion of all the atoms in an object. With this new understanding, it becomes possible to answer at one stroke the questions posed at the beginning of the section. Yes, it is at least theoretically possible to remove all the heat from an object. The coldest possible temperature, known as absolute zero, is that at which all the atoms have zero velocity, so that their kinetic energies, $K = (1/2)mv^2$, are all zero. No, there is no maximum amount of heat that a certain quantity of matter can have, and no maximum to the temperature scale, since arbitrarily large values of v can create arbitrarily large amounts of kinetic energy per atom.

The kinetic theory of heat also provides a simple explanation of the true nature of temperature. Temperature is a measure of the amount of energy per molecule, whereas heat is the total amount of energy possessed by all the molecules in an object.

There is an entire branch of physics, called thermodynamics, that deals with heat and temperature and forms the basis for technologies such as refrigeration. Thermodynamics is discussed in more detail in chapter 5, and I've provided here only a brief overview of the thermodynamic concepts that relate directly to energy.

2.4.2 All energy comes from particles moving or interacting.

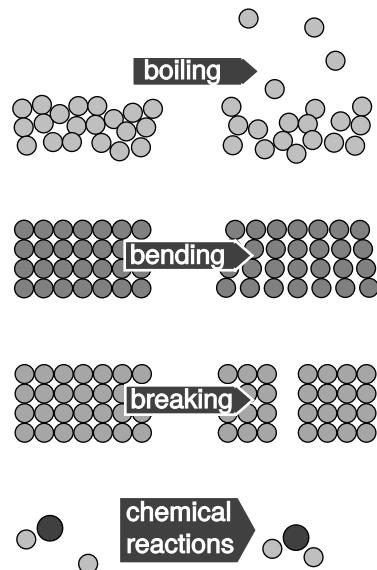
If I stretch the spring in figure c and then release it, it snaps taut again. The creation of some kinetic energy shows that there must have been some other form of energy that was destroyed. What was it?

We could just invent a new type of energy called “spring energy,” study its behavior, and call it quits, but that would be ugly. Are we going to have to invent a new forms of energy like this, over and over? No: the title of this book doesn’t lie, and physics really is fundamentally simple. As shown in figure d, when we bend or stretch an object, we’re really changing the distances between the atoms, resulting in a change in electrical energy. Electrical energy isn’t really our topic right now — that’s what most of the second half of this book is about — but conceptually it’s very similar to gravitational energy. Like gravitational energy, it depends on $1/r$, although there are some interesting new phenomena, such as the existence of both attraction and repulsion, which doesn’t occur with gravity because gravitational mass can’t be negative. The real point is that all the apparently dissimilar forms of energy in figure d turn out to be due to electrical interactions among atoms. Even if we wish to include nuclear reactions (figure e) in the picture, there still turn out to be only four fundamental types of energy:

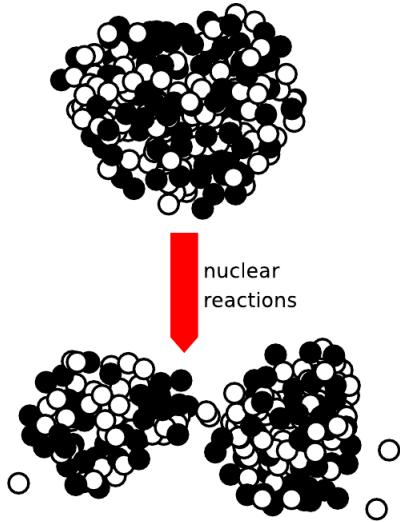
- kinetic energy** (including heat)
- gravitational energy**
- electrical and magnetic energy**
- nuclear energy**



c / The spring’s energy is really due to electrical interactions among atoms.



d / All these energy transformations turn out at the atomic level to be due to changes in the distances between atoms that interact electrically.



e / This figure looks similar to the previous ones, but the scale is a million times smaller. The little balls are the neutrons and protons that make up the tiny nucleus at the center of a uranium atom. When the nucleus splits (fissions), the source of the kinetic energy is partly electrical and partly nuclear.

Astute students have often asked me how light fits into this picture. This is a very good question, and in fact it could be argued that it is the basic question that led to Einstein's theory of relativity as well as the modern quantum picture of nature. Since these are topics for the second half of the book, we'll have to be content with half an answer at this point. For now, we may think of light energy as a form of kinetic energy, but one calculated not according to $(1/2)mv^2$ but by some other equation. (We know that $(1/2)mv^2$ would not make sense, because light has no mass, and furthermore, high-energy beams of light do not differ in speed from low-energy ones.)

Temperature during boiling

example 21

▷ If you stick a thermometer in a pan of water, and watch the temperature as you bring the water to a boil, you'll notice an interesting fact. The temperature goes up until the boiling point is reached, but then stays at 100°C during the whole time the water is being boiled off. The temperature of the steam is also 100°C. Why does the temperature "stick" like this? What's happening to all the energy that the stove's burner is putting into the pan?

▷ As shown in figure d, boiling requires an increase in electrical energy, because the atoms coming out as gas are moving away from the other atoms, which attract them electrically. It is only this electrical energy that is increasing, not the atoms' kinetic energy, which is what the thermometer can measure.

Diffusion

example 22

▷ A drop of food coloring in a cup of water will gradually spread out, even if you don't do any mixing with a spoon. This is called diffusion. Why would this happen, and what effect would temperature have? What would happen with solids or gases?

▷ Figure b shows that the atoms in a liquid mingle because of their random thermal motion. Diffusion is slow (typically on the order of a centimeter a minute), despite the *high* speeds of the atoms (typically hundreds of miles per hour). This is due to the randomness of the motion: a particular atom will take a long time to travel any significant distance, because it doesn't travel in a straight line.

Based on this picture, we expect that the speed of diffusion should increase as a function of temperature, and experiments show that this is true.

Diffusion also occurs in gases, which is why you can smell things even when the air is still. The speeds are much faster, because the typical distance between collisions is much longer than in a liquid.

We can see from figure b that diffusion won't occur in solids, because each atom vibrates around an equilibrium position.

Discussion Questions

A I'm not making this up. XS Energy Drink has ads that read like this: *All the "Energy" ... Without the Sugar! Only 8 Calories!* Comment on this.

2.4.3 Applications

Heat transfer

▷ Conduction

When you hold a hot potato in your hand, energy is transferred from the hot object to the cooler one. Our microscopic picture of this process (figure b, p. 110) tells us that the heat transfer can only occur at the surface of contact, where one layer of atoms in the potato skin make contact with one such layer in the hand. This type of heat transfer is called *conduction*, and its rate is proportional to both the surface area and the temperature difference.

▷ Convection

In a gas or a liquid, a faster method of heat transfer can occur, because hotter or colder parts of the fluid can flow, physically transporting their heat energy from one place to another. This mechanism of heat transfer, *convection*, is at work in Los Angeles when hot Santa Ana winds blow in from the Mojave Desert. On a cold day, the reason you feel warmer when there is no wind is that your skin warms a thin layer of air near it by conduction. If a gust of wind comes along, convection robs you of this layer. A thermos bottle has inner and outer walls separated by a layer of vacuum, which prevents heat transport by conduction or convection, except for a tiny amount of conduction through the thin connection between the walls, near the neck, which has a small cross-sectional area.

▷ Radiation

The glow of the sun or a candle flame is an example of heat transfer by *radiation*. In this context, “radiation” just means anything that radiates outward from a source, including, in these examples, ordinary visible light. The power is proportional to the surface area of the radiating object. It also depends very dramatically on the radiator’s absolute temperature, $P \propto T^4$.

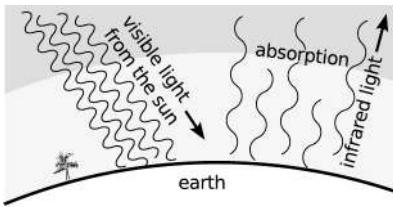
We can easily understand the reason for radiation based on the picture of heat as random kinetic energy at the atomic scale. Atoms are made out of subatomic particles, such as electrons and nuclei, that carry electric charge. When a charged particle vibrates, it creates wave disturbances in the electric and magnetic fields, and the waves have a frequency (number of vibrations per second) that matches the frequency of the particle’s motion. If this frequency is in the right range, they constitute visible light. When an object is closer to room temperature, it glows in the invisible infrared part of the spectrum.

Earth's energy equilibrium

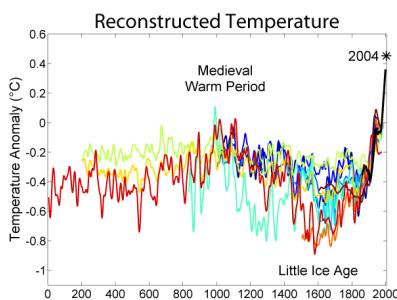
Our planet receives a nearly constant amount of energy from the sun (about 1.8×10^{17} W). If it hadn't had any mechanism for getting rid of that energy, the result would have been some kind of catastrophic explosion soon after its formation. Even a 10% imbalance between energy input and output, if maintained steadily from the time of the Roman Empire until the present, would have been enough to raise the oceans to a boil. So evidently the earth does dump this energy somehow. How does it do it? Our planet is surrounded by the vacuum of outer space, like the ultimate thermos bottle. Therefore it can't expel heat by conduction or convection, but it does radiate in the infrared, and this is the *only* available mechanism for cooling.

Global warming

It was realized starting around 1930 that this created a dangerous vulnerability in our biosphere. Our atmosphere is only about 0.04% carbon dioxide, but carbon dioxide is an extraordinarily efficient absorber of infrared light. It is, however, transparent to visible light. Therefore any increase in the concentration of carbon dioxide would decrease the efficiency of cooling by radiation, while allowing in just as much heat input from visible light. When we burn fossil fuels such as gasoline or coal, we release into the atmosphere carbon that had previously been locked away underground. This results in a shift to a new energy balance. The average temperature T of the land and oceans increases until the T^4 dependence of radiation compensates for the additional absorption of infrared light.



f / The "greenhouse effect." Carbon dioxide in the atmosphere allows visible light in, but partially blocks the reemitted infrared light.



g / Global average temperatures over the last 2000 years. The black line is from thermometer measurements. The colored lines are from various indirect indicators such as tree rings, ice cores, buried pollen, and corals.

By about 1980, a clear scientific consensus had emerged that this effect was real, that it was caused by human activity, and that it had resulted in an abrupt increase in the earth's average temperature. We know, for example, from radioisotope studies that the effect has not been caused by the release of carbon dioxide in volcanic eruptions. The temperature increase has been verified by multiple independent methods, including studies of tree rings and coral reefs. Detailed computer models have correctly predicted a number of effects that were later verified empirically, including a rise in sea levels, and day-night and pole-equator variations. There is no longer any controversy among climate scientists about the existence or cause of the effect.

One solution to the problem is to replace fossil fuels with renewable sources of energy such as solar power and wind. However, these cannot be brought online fast enough to prevent severe warming in the next few decades, so nuclear power is also a critical piece of the puzzle.

2.5 Oscillations

Let's revisit the example of the stretched spring from the previous section. We know that its energy is a form of electrical energy of interacting atoms, which is nice conceptually but doesn't help us to solve problems, since we don't know how the energy, U , depends on the length of the spring. All we know is that there's an equilibrium (figure a/1), which is a local minimum of the function U . An extremely important problem which arises in this connection is how to calculate oscillatory motion around an equilibrium, as in a/4-13. Even if we did special experiments to find out how the spring's energy worked, it might seem like we'd have to go through just as much work to deal with any other kind of oscillation, such as a sapling swinging back and forth in the breeze.

Surprisingly, it's possible to analyze this type of oscillation in a very general and elegant manner, as long as the analysis is limited to *small* oscillations. We'll talk about the mass on the spring for concreteness, but there will be nothing in the discussion at all that is restricted to that particular physical system. First, let's choose a coordinate system in which $x = 0$ corresponds to the position of the mass where the spring is in equilibrium, and since interaction energies like U are only well defined up to an additive constant, we'll simply define it to be zero at equilibrium:

$$U(0) = 0$$

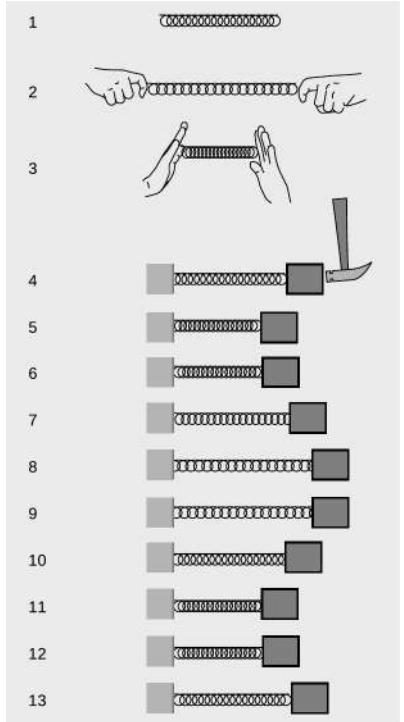
Since $x = 0$ is an equilibrium, $U(x)$ must have a local minimum there, and a differentiable function (which we assume U is) has a zero derivative at a local minimum:

$$\frac{dU}{dx}(0) = 0$$

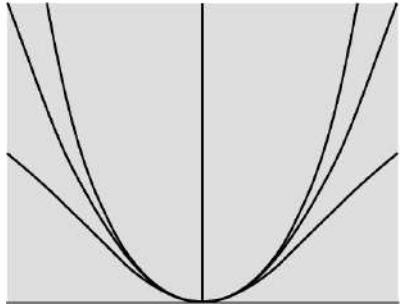
There are still infinitely many functions that could satisfy these criteria, including the three shown in figure b, which are $x^2/2$, $x^2/2(1+x^2)$, and $(e^{3x} + e^{-3x} - 2)/18$. Note, however, how all three functions are virtually identical right near the minimum. That's because they all have the same curvature. More specifically, each function has its second derivative equal to 1 at $x = 0$, and the second derivative is a measure of curvature. We write k for the second derivative of the energy at an equilibrium point,

$$\frac{d^2U}{dx^2}(0) = k.$$

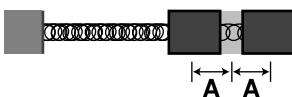
Physically, k is a measure of stiffness. For example, the heavy-duty springs in a car's shock absorbers would have a high value of k . It is often referred to as the spring constant, but we're only using a spring as an example here. As shown in figure b, any two functions that have $U(0) = 0$, $dU/dx = 0$, and $d^2U/dx^2 = k$, with the same value of k , are virtually indistinguishable for small values of x , so if



a / The spring has a minimum-energy length, 1, and energy is required in order to compress or stretch it, 2 and 3. A mass attached to the spring will oscillate around the equilibrium, 4-13.



b / Three functions with the same curvature at $x=0$.



c / The amplitude would usually be defined as the distance from equilibrium to one extreme of the motion, i.e., half the total travel.

we want to analyze small oscillations, it doesn't even matter which function we assume. For simplicity, we'll just use $U(x) = (1/2)kx^2$ from now on.

Now we're ready to analyze the mass-on-a-spring system, while keeping in mind that it's really only a representative example of a whole class of similar oscillating systems. We expect that the motion is going to repeat itself over and over again, and since we're not going to include frictional heating in our model, that repetition should go on forever without dying out. The most interesting thing to know about the motion would be the period, T , which is the amount of time required for one complete cycle of the motion. We might expect that the period would depend on the spring constant, k , the mass, m , and the amplitude, A , defined in figure c.¹¹

In examples like the brachistochrone and the Apollo 11 mission, it was generally necessary to use numerical techniques to determine the amount of time required for a certain motion. Once again, let's dust off the `time3` function from page 93 and modify it for our purposes. For flexibility, we'll define the function $U(x)$ as a separate Python function. We really want to calculate the time required for the mass to come back to its starting point, but that would be awkward to set up, since our function works by dividing up the distance to be traveled into tiny segments. By symmetry, the time required to go from one end to the other equals the time required to come back to the start, so we'll just calculate the time for half a cycle and then double it when we return the result at the end of the function. The test at lines 16-19 is necessary because otherwise at the very end of the motion we can end up trying to take the square root of a negative number due to rounding errors.

¹¹Many kinds of oscillations are possible, so there is no standard definition of the amplitude. For a pendulum, the natural definition would be in terms of an angle. For a radio transmitter, we'd use some kind of electrical units.

```

1 import math
2 def u(k,x):
3     return .5*k*x**2
4
5 def osc(m,k,a,n):
6     x=a
7     v=0
8     dx = -2.*a/n
9     t=0
10    e = u(k,x)+.5*m*v**2
11    for i in range(n):
12        x_old = x
13        v_old = v
14        x = x+dx
15        kinetic = e-u(k,x)
16        if kinetic<0. :
17            v=0.
18            print "warning, K=",kinetic,"<0"
19        else :
20            v = -math.sqrt(2.*kinetic/m)
21            v_avg = (v+v_old)/2.
22            dt=dx/v_avg
23            t=t+dt
24    return 2.*t
25

>>> print(osc(1.,1.,1.,100000))
warning, K= -1.43707268307e-12 <0
6.2831854132667919

```

The first thing to notice is that with this particular set of inputs ($m=1$ kg, $k = 1$ J/m², and $A = 1$ m), the program has done an excellent job of computing $2\pi = 6.2831853\dots$. This is Mother Nature giving us a strong hint that the problem has an algebraic solution, not just a numerical one. The next interesting thing happens when we change the amplitude from 1 m to 2 m:

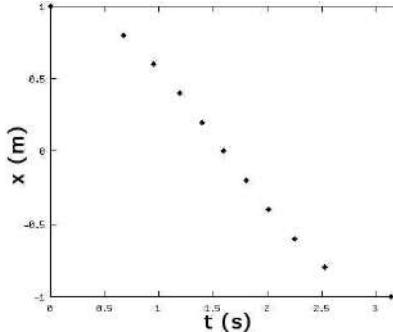
```

>>> print(osc(1.,1.,2.,100000))
warning, K= -5.7482907323e-12 <0
6.2831854132667919

```

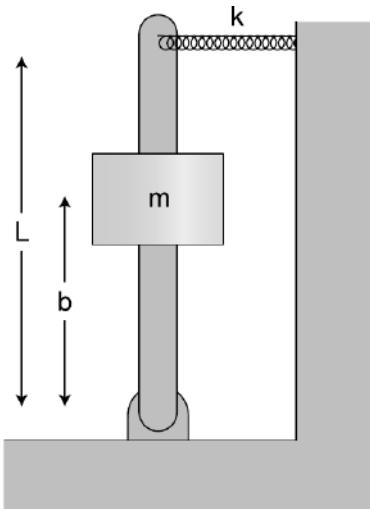
Even though the mass had to travel double the distance in each direction, the period is the same to within the numerical accuracy of the calculation!

With these hints, it seems like we should start looking for an algebraic solution. For guidance, here's a graph of x as a function of t , as calculated by the `osc` function with $n=10$.



This looks like a cosine function, so let's see if a $x = A \cos(\omega t + \delta)$ is a solution to the conservation of energy equation — it's not uncommon to try to "reverse-engineer" the cryptic results of a numerical calculation like this. The symbol $\omega = 2\pi/T$ (Greek omega), called angular frequency, is a standard symbol for the number of radians per second of oscillation. Except for the factor of 2π , it is identical to the ordinary frequency $f = 1/T$, which has units of s^{-1} or Hz (Hertz). The phase angle δ is to allow for the possibility that $t = 0$ doesn't coincide with the beginning of the motion. The energy is

$$\begin{aligned} E &= K + U \\ &= \frac{1}{2}mv^2 + \frac{1}{2}kx^2 \\ &= \frac{1}{2}m\left(\frac{dx}{dt}\right)^2 + \frac{1}{2}kx^2 \\ &= \frac{1}{2}m[-A\omega \sin(\omega t + \delta)]^2 + \frac{1}{2}k[A \cos(\omega t + \delta)]^2 \\ &= \frac{1}{2}A^2[m\omega^2 \sin^2(\omega t + \delta) + k \cos^2(\omega t + \delta)] \end{aligned}$$



d / Example 23. The rod pivots on the hinge at the bottom.

According to conservation of energy, this has to be a constant. Using the identity $\sin^2 + \cos^2 = 1$, we can see that it will be a constant if we have $m\omega^2 = k$, or $\omega = \sqrt{k/m}$, i.e., $T = 2\pi\sqrt{m/k}$. Note that the period is independent of amplitude.

A spring and a lever

example 23

▷ What is the period of small oscillations of the system shown in the figure? Neglect the mass of the lever and the spring. Assume that the spring is so stiff that gravity is not an important effect. The spring is relaxed when the lever is vertical.

▷ This is a little tricky, because the spring constant k , although it is relevant, is *not* the k we should be putting into the equation $T = 2\pi\sqrt{m/k}$. The k that goes in there has to be the second derivative of U with respect to the position, x , of the mass that's moving. The energy U stored in the spring depends on how far the *tip* of the lever is from the center. This distance equals $(L/b)x$,

so the energy in the spring is

$$U = \frac{1}{2}k \left(\frac{L}{b}x\right)^2 \\ = \frac{kL^2}{2b^2}x^2,$$

and the k we have to put in $T = 2\pi\sqrt{m/k}$ is

$$\frac{d^2 U}{dx^2} = \frac{kL^2}{b^2}.$$

The result is

$$T = 2\pi\sqrt{\frac{mb^2}{kL^2}} \\ = \frac{2\pi b}{L}\sqrt{\frac{m}{k}}$$

The leverage of the lever makes it as if the spring was stronger, and decreases the period of the oscillations by a factor of b/L .

Water in a U-shaped tube

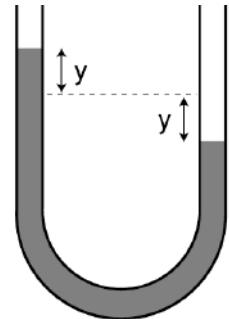
example 24

- ▷ What is the period of oscillation of the water in figure e?
- ▷ In example 13 on p. 89, we found $U(y) = \rho g A y^2$, so the “spring constant,” which really isn’t a spring constant here at all, is

$$k = \frac{d^2 U}{dy^2} \\ = 2\rho g A.$$

This is an interesting example, because k can be calculated without any approximations, but the kinetic energy requires an approximation, because we don’t know the details of the pattern of flow of the water. It could be very complicated. There will be a tendency for the water near the walls to flow more slowly due to friction, and there may also be swirling, turbulent motion. However, if we make the approximation that all the water moves with the same velocity as the surface, dy/dt , then the mass-on-a-spring analysis applies. Letting L be the total length of the filled part of the tube, the mass is $\rho L A$, and we have

$$T = 2\pi\sqrt{m/k} \\ = 2\pi\sqrt{\frac{\rho L A}{2\rho g A}} \\ = 2\pi\sqrt{\frac{L}{2g}}.$$



e / Water in a U-shaped tube.

This chapter is summarized on page 1076. Notation and terminology are tabulated on pages 1070-1071.

Problems

The symbols \checkmark , \blacksquare , etc. are explained on page 127.

1 Experiments show that the power consumed by a boat's engine is approximately proportional to the third power of its speed. (We assume that it is moving at constant speed.)

(a) When a boat is cruising at constant speed, what type of energy transformation do you think is being performed?

(b) If you upgrade to a motor with double the power, by what factor is your boat's maximum cruising speed increased?

► Solution, p. 1039 \blacksquare

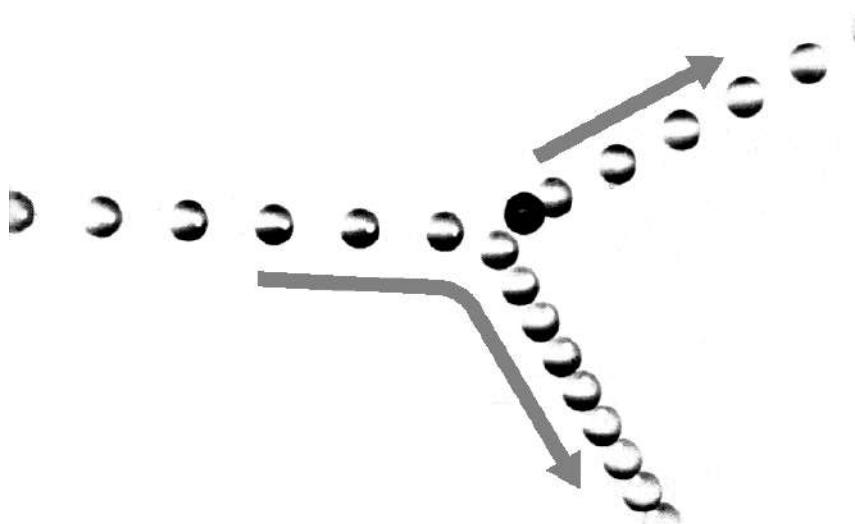
2 Object A has a kinetic energy of 13.4 J. Object B has a mass that is greater by a factor of 3.77, but is moving more slowly by a factor of 2.34. What is object B's kinetic energy?

► Solution, p. 1039 \blacksquare

3 My 1.25 kW microwave oven takes 126 seconds to bring 250 g of water from room temperature to a boil. What percentage of the power is being wasted? Where might the rest of the energy be going?

► Solution, p. 1039 \blacksquare

4 The multiflash photograph shows a collision between two pool balls. The ball that was initially at rest shows up as a dark image in its initial position, because its image was exposed several times before it was struck and began moving. By making *measurements* on the figure, determine *numerically* whether or not energy appears to have been conserved in the collision. What systematic effects would limit the accuracy of your test? [From an example in PSSC Physics.]



Problem 4.

5 A grasshopper with a mass of 110 mg falls from rest from a height of 310 cm. On the way down, it dissipates 1.1 mJ of heat due to air resistance. At what speed, in m/s, does it hit the ground?

► Solution, p. 1039 ■

6 A ball rolls up a ramp, turns around, and comes back down. When does it have the greatest gravitational energy? The greatest kinetic energy? [Based on a problem by Serway and Faughn.] ■

7 (a) You release a magnet on a tabletop near a big piece of iron, and the magnet leaps across the table to the iron. Does the magnetic energy increase, or decrease? Explain. (b) Suppose instead that you have two repelling magnets. You give them an initial push towards each other, so they decelerate while approaching each other. Does the magnetic energy increase, or decrease? Explain. ■

8 Estimate the kinetic energy of an Olympic sprinter. ■

9 You are driving your car, and you hit a brick wall head on, at full speed. The car has a mass of 1500 kg. The kinetic energy released is a measure of how much destruction will be done to the car and to your body. Calculate the energy released if you are traveling at (a) 40 mi/hr, and again (b) if you're going 80 mi/hr. What is counterintuitive about this, and what implication does this have for driving at high speeds? ✓ ■

10 A closed system can be a bad thing — for an astronaut sealed inside a space suit, getting rid of body heat can be difficult. Suppose a 60-kg astronaut is performing vigorous physical activity, expending 200 W of power. If none of the heat can escape from her space suit, how long will it take before her body temperature rises by 6°C (11°F), an amount sufficient to kill her? Assume that the amount of heat required to raise her body temperature by 1°C is the same as it would be for an equal mass of water. Express your answer in units of minutes. ✓ ■

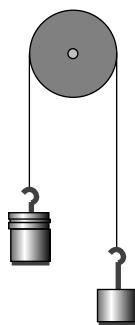
11 The following table gives the amount of energy required in order to heat, melt, or boil a gram of water.

heat 1 g of ice by 1°C	2.05 J
melt 1 g of ice	333 J
heat 1 g of water by 1°C	4.19 J
boil 1 g of water	2500 J
heat 1 g of steam by 1°C	2.01 J

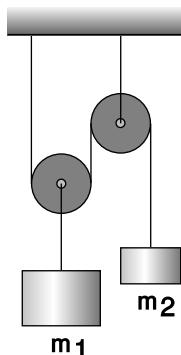
(a) How much energy is required in order to convert 1.00 g of ice at -20°C into steam at 137°C ? ✓

(b) What is the minimum amount of hot water that could melt 1.00 g of ice? ✓ ■

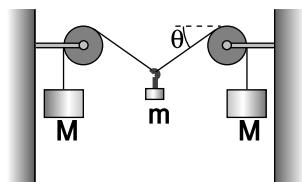
12 Anya climbs to the top of a tree, while Ivan climbs half-way to the top. They both drop pennies to the ground. Compare the kinetic energies and velocities of the pennies on impact, using ratios. ■



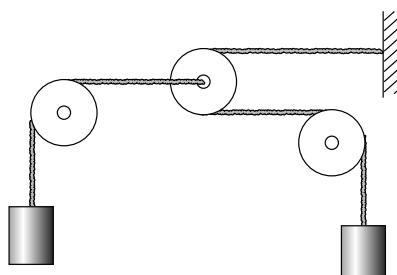
Problem 16.



Problem 17.



Problem 18.



Problem 19.

13 Anya and Ivan lean over a balcony side by side. Anya throws a penny downward with an initial speed of 5 m/s. Ivan throws a penny upward with the same speed. Both pennies end up on the ground below. Compare their kinetic energies and velocities on impact. ■

14 (a) A circular hoop of mass m and radius r spins like a wheel while its center remains at rest. Let ω (Greek letter omega) be the number of radians it covers per unit time, i.e., $\omega = 2\pi/T$, where the period, T , is the time for one revolution. Show that its kinetic energy equals $(1/2)m\omega^2r^2$.

(b) Show that the answer to part a has the right units. (Note that radians aren't really units, since the definition of a radian is a unitless ratio of two lengths.)

(c) If such a hoop rolls with its center moving at velocity v , its kinetic energy equals $(1/2)mv^2$, plus the amount of kinetic energy found in part a. Show that a hoop rolls down an inclined plane with half the acceleration that a frictionless sliding block would have. ■

15 On page 83, I used the chain rule to prove that the acceleration of a free-falling object is given by $a = -g$. In this problem, you'll use a different technique to prove the same thing. Assume that the acceleration is a constant, a , and then integrate to find v and y , including appropriate constants of integration. Plug your expressions for v and y into the equation for the total energy, and show that $a = -g$ is the only value that results in constant energy. ■

16 The figure shows two unequal masses, m_1 and m_2 , connected by a string running over a pulley. Find the acceleration.

▷ Hint, p. 1034 ✓ ■

17

What ratio of masses will balance the pulley system shown in the figure?

▷ Hint, p. 1034 ■

18 (a) For the apparatus shown in the figure, find the equilibrium angle θ in terms of the two masses. ✓

(b) Interpret your result in the case of $M \gg m$ (M much greater than m). Does it make sense physically?

(c) For what combinations of masses would your result give nonsense? Interpret this physically. ▷ Hint, p. 1034 ■

19 In the system shown in the figure, the pulleys on the left and right are fixed, but the pulley in the center can move to the left or right. The two hanging masses are identical, and the pulleys and ropes are all massless. Find the upward acceleration of the mass on the left, in terms of g only. ▷ Hint, p. 1034 ✓ ■

20 Two atoms will interact through electrical forces between their protons and electrons. One fairly good approximation to the electrical energy is the Lennard-Jones formula,

$$U(r) = k \left[\left(\frac{a}{r} \right)^{12} - 2 \left(\frac{a}{r} \right)^6 \right],$$

where r is the center-to-center distance between the atoms and k is a positive constant. Show that (a) there is an equilibrium point at $r = a$,

- (b) the equilibrium is stable, and
- (c) the energy required to bring the atoms from their equilibrium separation to infinity is k .

▷ Hint, p. 1034 ■

21 The International Space Station orbits at an altitude of about 360 to 400 km. What is the gravitational field of the earth at this altitude? ✓ ■

22 (a) A geosynchronous orbit is one in which the satellite orbits above the equator, and has an orbital period of 24 hours, so that it is always above the same point on the spinning earth. Calculate the altitude of such a satellite. ✓

(b) What is the gravitational field experienced by the satellite? Give your answer as a percentage in relation to the gravitational field at the earth's surface.

▷ Hint, p. 1034 ✓ ■

23 Astronomers calculating orbits of planets often work in a nonmetric system of units, in which the unit of time is the year, the unit of mass is the sun's mass, and the unit of distance is the astronomical unit (A.U.), defined as half the long axis of the earth's orbit. In these units, find an exact expression for the gravitational constant, G . ✓ ■

24 The star Lalande 21185 was found in 1996 to have two planets in roughly circular orbits, with periods of 6 and 30 years. What is the ratio of the two planets' orbital radii? ✓ ■

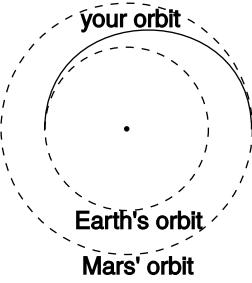
25 A projectile is moving directly away from a planet of mass M at exactly escape velocity. (a) Find r , the distance from the projectile to the center of the planet, as a function of time, t , and also find $v(t)$. ✓

(b) Check the units of your answer.

(c) Does v show the correct behavior as t approaches infinity?

▷ Hint, p. 1034 ■

26 The purpose of this problem is to estimate the height of the tides. The main reason for the tides is the moon's gravity, and we'll neglect the effect of the sun. Also, real tides are heavily influenced by landforms that channel the flow of water, but we'll think of the earth as if it was completely covered with oceans. Under these assumptions, the ocean surface should be a surface of constant U/m . That is, a thimbleful of water, m , should not be able to gain or lose



Problem 27.

any gravitational energy by moving from one point on the ocean surface to another. If only the spherical earth's gravity was present, then we'd have $U/m = -GM_e/r$, and a surface of constant U/m would be a surface of constant r , i.e., the ocean's surface would be spherical. Taking into account the moon's gravity, the main effect is to shift the center of the sphere, but the sphere also becomes slightly distorted into an approximately ellipsoidal shape. (The shift of the center is not physically related to the tides, since the solid part of the earth tends to be centered within the oceans; really, this effect has to do with the motion of the whole earth through space, and the way that it wobbles due to the moon's gravity.) Determine the amount by which the long axis of the ellipsoid exceeds the short axis. ▷ Hint, p. 1034 ■

27 You are considering going on a space voyage to Mars, in which your route would be half an ellipse, tangent to the Earth's orbit at one end and tangent to Mars' orbit at the other. Your spacecraft's engines will only be used at the beginning and end, not during the voyage. How long would the outward leg of your trip last? (The orbits of Earth and Mars are nearly circular, and Mars's is bigger by a factor of 1.52.) ✓ ■

28 When you buy a helium-filled balloon, the seller has to inflate it from a large metal cylinder of the compressed gas. The helium inside the cylinder has energy, as can be demonstrated for example by releasing a little of it into the air: you hear a hissing sound, and that sound energy must have come from somewhere. The total amount of energy in the cylinder is very large, and if the valve is inadvertently damaged or broken off, the cylinder can behave like a bomb or a rocket.

Suppose the company that puts the gas in the cylinders prepares cylinder A with half the normal amount of pure helium, and cylinder B with the normal amount. Cylinder B has twice as much energy, and yet the temperatures of both cylinders are the same. Explain, at the atomic level, what form of energy is involved, and why cylinder B has twice as much. ■

29 Explain in terms of conservation of energy why sweating cools your body, even though the sweat is at the same temperature as your body. Describe the forms of energy involved in this energy transformation. Why don't you get the same cooling effect if you wipe the sweat off with a towel? Hint: The sweat is evaporating. ■

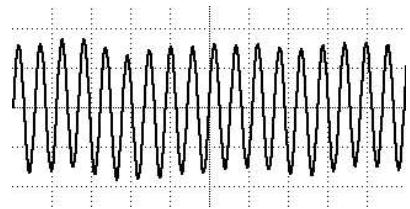
30

[This problem is now problem 3-73.] ■

31 All stars, including our sun, show variations in their light

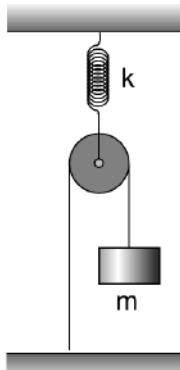
output to some degree. Some stars vary their brightness by a factor of two or even more, but our sun has remained relatively steady during the hundred years or so that accurate data have been collected. Nevertheless, it is possible that climate variations such as ice ages are related to long-term irregularities in the sun's light output. If the sun was to increase its light output even slightly, it could melt enough Antarctic ice to flood all the world's coastal cities. The total sunlight that falls on Antarctica amounts to about 1×10^{16} watts. In the absence of natural or human-caused climate change, this heat input to the poles is balanced by the loss of heat via winds, ocean currents, and emission of infrared light, so that there is no net melting or freezing of ice at the poles from year to year. Suppose that the sun changes its light output by some small percentage, but there is no change in the rate of heat loss by the polar caps. Estimate the percentage by which the sun's light output would have to increase in order to melt enough ice to raise the level of the oceans by 10 meters over a period of 10 years. (This would be enough to flood New York, London, and many other cities.) Melting 1 kg of ice requires 3×10^3 J.

32 The figure shows the oscillation of a microphone in response to the author whistling the musical note "A." The horizontal axis, representing time, has a scale of 1.0 ms per square. Find the period T , the frequency f , and the angular frequency ω . \checkmark



Problem 32.

33 (a) A mass m is hung from a spring whose spring constant is k . Write down an expression for the total interaction energy of the system, U , and find its equilibrium position. \triangleright Hint, p. 1034 \checkmark
 (b) Explain how you could use your result from part a to determine an unknown spring constant. \square

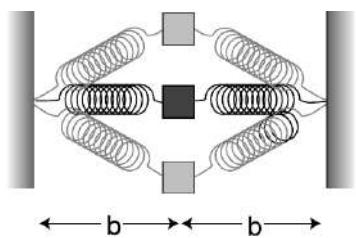


Problem 35.

34 A certain mass, when hung from a certain spring, causes the spring to stretch by an amount h compared to its equilibrium length. If the mass is displaced vertically from this equilibrium, it will oscillate up and down with a period T_{osc} . Give a numerical comparison between T_{osc} and T_{fall} , the time required for the mass to fall from rest through a height h , when it isn't attached to the spring. (You will need the result of problem 33). \checkmark

35 Find the period of vertical oscillations of the mass m . The spring, pulley, and ropes have negligible mass.

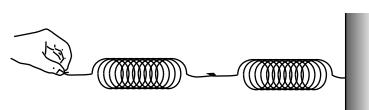
\triangleright Hint, p. 1034 \checkmark



Problem 36.

36 The equilibrium length of each spring in the figure is b , so when the mass m is at the center, neither spring exerts any force on it. When the mass is displaced to the side, the springs stretch; their spring constants are both k .

- Find the energy, U , stored in the springs, as a function of y , the distance of the mass up or down from the center. \checkmark
- Show that the period of small up-down oscillations is infinite. \square



Problem 37.

37 Two springs with spring constants k_1 and k_2 are put together end-to-end. Let x_1 be the amount by which the first spring is stretched relative to its equilibrium length, and similarly for x_2 . If the combined double spring is stretched by an amount b relative to its equilibrium length, then $b = x_1 + x_2$. Find the spring constant, K , of the combined spring in terms of k_1 and k_2 .

▷ Hint, p. 1035 ▷ Answer, p. 1068 ✓

38 A mass m on a spring oscillates around an equilibrium at $x = 0$. Any function $U(x)$ with an equilibrium at $x = 0$ can be approximated as $U(x) = (1/2)kx^2$, and if the energy is symmetric with respect to positive and negative values of x , then the next level of improvement in such an approximation would be $U(x) = (1/2)kx^2 + bx^4$. The general idea here is that any smooth function can be approximated locally by a polynomial, and if you want a better approximation, you can use a polynomial with more terms in it. When you ask your calculator to calculate a function like \sin or e^x , it's using a polynomial approximation with 10 or 12 terms. Physically, a spring with a positive value of b gets stiffer when stretched strongly than an "ideal" spring with $b = 0$. A spring with a negative b is like a person who cracks under stress — when you stretch it too much, it becomes more elastic than an ideal spring would. We should not expect any spring to give totally ideal behavior no matter how much it is stretched. For example, there has to be some point at which it breaks.

Do a numerical simulation of the oscillation of a mass on a spring whose energy has a nonvanishing b . Is the period still independent of amplitude? Is the amplitude-independent equation for the period still approximately valid for small enough amplitudes? Does the addition of a positive x^4 term tend to increase the period, or decrease it? Include a printout of your program and its output with your homework paper. ■

39 An idealized pendulum consists of a pointlike mass m on the end of a massless, rigid rod of length L . Its amplitude, θ , is the angle the rod makes with the vertical when the pendulum is at the end of its swing. Write a numerical simulation to determine the period of the pendulum for any combination of m , L , and θ . Examine the effect of changing each variable while manipulating the others. ■

40 A ball falls from a height h . Without air resistance, the time it takes to reach the floor is $\sqrt{2h/g}$. A numerical version of this calculation was given in program `time2` on page 92. Now suppose that air resistance is not negligible. For a smooth sphere of radius r , moving at speed v through air of density ρ , the amount of energy dQ dissipated as heat as the ball falls through a height dy is given (ignoring signs) by $dQ = (\pi/4)\rho v^2 r^2 dy$. Modify the program to incorporate this effect, and find the resulting change in the fall

time in the case of a 21 g ball of radius 1.0 cm, falling from a height of 1.0 m. The density of air at sea level is about 1.2 kg/m^3 . Turn in a printout of both your program and its output. Answer: 0.34 ms.



41 The factorial of an integer n , written $n!$, is defined as the product of all the positive integers less than or equal to n . For example, $3! = 1 \times 2 \times 3 = 6$. Write a Python program to compute the factorial of a number. Test it with a small number whose factorial you can check by hand. Then use it to compute $30!$. (Python computes integer results with unlimited precision, so you won't get any problems with rounding or overflows.) Turn in a printout of your program and its output, including the test.



42 Estimate the kinetic energy of a buzzing fly's wing. (You may wish to review subsection 0.2.3 on order-of-magnitude estimates.)



43 A blade of grass moves upward as it grows. Estimate its kinetic energy. (You may wish to review subsection 0.2.3 on order-of-magnitude estimates.)



Key to symbols:

■ easy ■ typical ■ challenging ■ difficult ■ very difficult

✓ An answer check is available at www.lightandmatter.com.

Exercises

Exercise 2A: Reasoning with Ratios and Powers

Equipment:

ping-pong balls and paddles

two-meter sticks

You have probably bounced a ping pong ball straight up and down in the air. The time between hits is related to the height to which you hit the ball. If you take twice as much time between hits, how many times higher do you think you will have to hit the ball? Write down your hypothesis:_____

Your instructor will first beat out a tempo of 240 beats per minute (four beats per second), which you should try to match with the ping-pong ball. Measure the height to which the ball rises:_____

Now try it at 120 beats per minute:_____

Compare your hypothesis and your results with the rest of the class.

Exercise 2B: The Shell Theorem

This exercise is an approximate numerical test of the shell theorem. There are seven masses A-G, each being one kilogram. Masses A-E, each one meter from the center, form a shape like two Egyptian pyramids joined at their bases; this is a rough approximation to a six-kilogram spherical shell of mass. Mass G is five meters from the center of the main group. The class will divide into six groups and split up the work required in order to calculate the total gravitational energy of mass G.

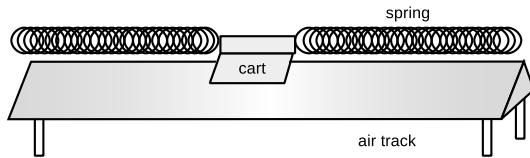


1. Each group should write its results on the board in units of picojoules, retaining six significant figures of precision.
2. The class will add the results and compare with the result that would be obtained with the shell theorem.

Exercise 2C: Vibrations

Equipment:

- air track and carts of two different masses
- springs
- weights



Place the cart on the air track and attach springs so that it can vibrate.

1. Test whether the period of vibration depends on amplitude. Try at least one moderate amplitude, for which the springs do not go slack, at least one amplitude that is large enough so that they do go slack, and one amplitude that's the very smallest you can possibly observe.
2. Try a cart with a different mass. Does the period change by the expected factor, based on the equation $T = 2\pi\sqrt{m/k}$?
3. In homework problem 33 on page 125, you showed that a spring's spring constant can be determined by hanging a weight from it. Use this technique to find the spring constant of each of the two springs. The equivalent spring constant of these two springs, attached to the cart in this way, can be found by adding their spring constants.
4. Test the equation $T = 2\pi\sqrt{m/k}$ numerically.



Forces transfer momentum to the girl.

Chapter 3

Conservation of Momentum

I think, therefore I am.

I hope that posterity will judge me kindly, not only as to the things which I have explained, but also to those which I have intentionally omitted so as to leave to others the pleasure of discovery.

René Descartes



a / Systems consisting of material particles that interact through an energy $U(r)$. *Top:* The galaxy M100. Here the “particles” are stars. *Middle:* The pool balls don’t interact until they come together and become compressed; the energy $U(r)$ has a sharp upturn when the center-to-center distance r gets small enough for the balls to be in contact. *Bottom:* A uranium nucleus undergoing fission. The energy $U(r)$ has a repulsive contribution from the electrical interactions of the protons, plus an attractive one due to the strong nuclear interaction. (*M100: Hubble Space Telescope image.*)

3.1 Momentum in one dimension

3.1.1 Mechanical momentum

In the martial arts movie *Crouching Tiger, Hidden Dragon*, those who had received mystical enlightenment are able to violate the laws of physics. Some of the violations are obvious, such as their ability to fly, but others are a little more subtle. The rebellious young heroine/antiheroine Jen Yu gets into an argument while sitting at a table in a restaurant. A young tough, Iron Arm Lu, comes running toward her at full speed, and she puts up one arm and effortlessly makes him bounce back, without even getting out of her seat or bracing herself against anything. She does all this between bites.

Although kinetic energy doesn’t depend on the direction of motion, we’ve already seen on page 89 how conservation of energy combined with Galilean relativity allows us to make some predictions about the direction of motion. One of the examples was a demonstration that it isn’t possible for a hockey puck to spontaneously reverse its direction of motion. In the scene from the movie, however, the woman’s assailant isn’t just gliding through space. He’s interacting with her, so the previous argument doesn’t apply here, and we need to generalize it to more than one object. We consider the case of a physical system composed of pointlike material particles, in which every particle interacts with every other particle through an energy $U(r)$ that depends only on the distance r between them. This still allows for a fairly general *mechanical system*, by which I mean roughly a system made of matter, not light. The characters in the movie are made of protons, neutrons, and electrons, so they would constitute such a system if the interactions among all these particles were of the form $U(r)$.¹ We might even be able to get away with thinking of each person as one big particle, if it’s a good approximation to say that every part of each person’s whole body moves in the same direction at the same speed.

The basic insight can be extracted from the special case where there are only two particles interacting, and they only move in one dimension, as in the example shown in figure b. Conservation of energy says

$$K_{1i} + K_{2i} + U_i = K_{1f} + K_{2f} + U_f.$$

For simplicity, let’s assume that the interactions start after the time we’re calling initial, and end before the instant we choose as final. This is true in figure b, for example. Then $U_i = U_f$, and we can subtract the interaction energies from both sides, giving,

$$\begin{aligned} K_{1i} + K_{2i} &= K_{1f} + K_{2f} \\ \frac{1}{2}m_1v_{1i}^2 + \frac{1}{2}m_2v_{2i}^2 &= \frac{1}{2}m_1v_{1f}^2 + \frac{1}{2}m_2v_{2f}^2. \end{aligned}$$

¹Electrical and magnetic interactions *don’t* quite behave like this, which is a point we’ll take up later in the book.

As in the one-particle argument on page 89, the trick is to require conservation of energy not just in one particular frame of reference, but in every frame of reference. In a frame of reference moving at velocity u relative to the first one, the velocities all have u added onto them:²

$$\frac{1}{2}m_1(v_{1i}+u)^2 + \frac{1}{2}m_2(v_{2i}+u)^2 = \frac{1}{2}m_1(v_{1f}+u)^2 + \frac{1}{2}m_2(v_{2f}+u)^2$$

When we square a quantity like $(v_{1i}+u)^2$, we get the same v_{1i}^2 that occurred in the original frame of reference, plus two u -dependent terms, $2v_{1i}u + u^2$. Subtracting the original conservation of energy equation from the version in the new frame of reference, we have

$$m_1v_{1i}u + m_2v_{2i}u = m_1v_{1f}u + m_2v_{2f}u,$$

or, dividing by u ,

$$m_1v_{1i} + m_2v_{2i} = m_1v_{1f} + m_2v_{2f}.$$

This is a statement that when you add up mv for the whole system, that total remains constant over time. In other words, this is a conservation law. The quantity mv is called *momentum*, notated p for obscure historical reasons. Its units are kg · m/s.

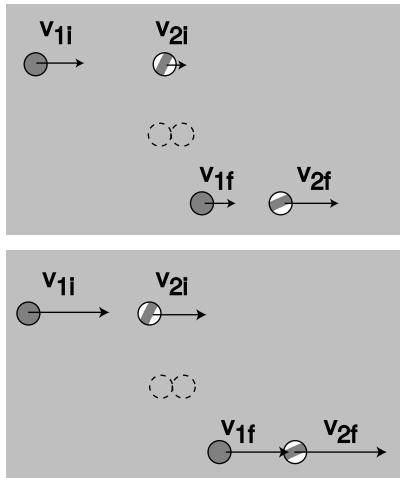
Unlike kinetic energy, momentum depends on the direction of motion, since the velocity is not squared. In one dimension, motion in the same direction as the positive x axis is represented with positive values of v and p . Motion in the opposite direction has negative v and p .

Jen Yu meets Iron Arm Lu

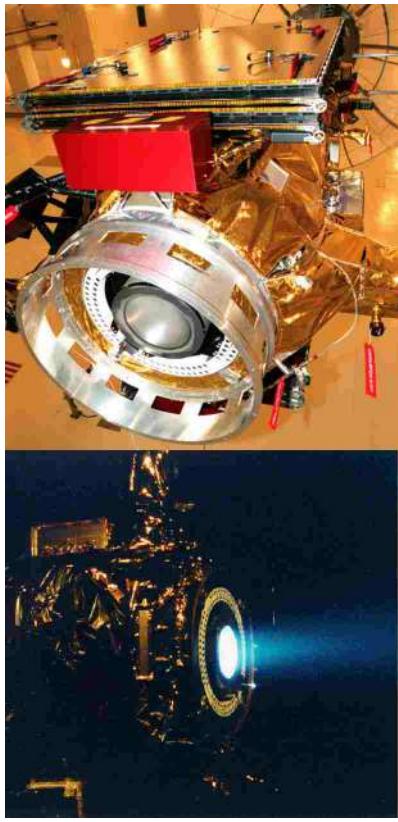
example 1

- ▷ Initially, Jen Yu is at rest, and Iron Arm Lu is charging to the left, toward her, at 5 m/s. Jen Yu's mass is 50 kg, and Lu's is 100 kg. After the collision, the movie shows Jen Yu still at rest, and Lu rebounding at 5 m/s to the right. Is this consistent with the laws of physics, or would it be impossible in real life?
- ▷ This is perfectly consistent with conservation of mass ($50\text{ kg}+100\text{ kg}=50\text{ kg}+100\text{ kg}$), and also with conservation of energy, since neither person's kinetic energy changes, and there is therefore no change in the total energy. (We don't have to worry about interaction energies, because the two points in time we're considering are ones at which the two people aren't interacting.) To analyze whether the scene violates conservation of momentum, we have to pick a coordinate system. Let's define positive as

²We can now see that the derivation would have been equally valid for $U_i \neq U_f$. The two observers agree on the distance between the particles, so they also agree on the interaction energies, even though they disagree on the kinetic energies.



b / A collision between two pool balls is seen in two different frames of reference. The solid ball catches up with the striped ball. Velocities are shown with arrows. The second observer is moving to the left at velocity u compared to the first observer, so all the velocities in the second frame have u added onto them. The two observers must agree on conservation of energy.



c / The ion drive engine of the NASA Deep Space 1 probe, shown under construction (top) and being tested in a vacuum chamber (bottom) prior to its October 1998 launch. Intended mainly as a test vehicle for new technologies, the craft nevertheless also carried out a scientific program that included a rendezvous with a comet in 2004. (NASA)

being to the right. The initial momentum is $(50 \text{ kg})(0 \text{ m/s}) + (100 \text{ kg})(-5 \text{ m/s}) = -500 \text{ kg} \cdot \text{m/s}$, and the final momentum is $(50 \text{ kg})(0 \text{ m/s}) + (100 \text{ kg})(5 \text{ m/s}) = 500 \text{ kg} \cdot \text{m/s}$. This is a change of 1000 $\text{kg} \cdot \text{m/s}$, which is impossible if the two people constitute a closed system.

One could argue that they're not a closed system, since Lu might be exchanging momentum with the floor, and Jen Yu might be exchanging momentum with the seat of her chair. This is a reasonable objection, but in the following section we'll see that there are physical reasons why, in this situation, the force of friction would be relatively weak, and would not be able to transfer that much momentum in a fraction of a second.

This example points to an intuitive interpretation of conservation of momentum, which is that interactions are always mutual. That is, Jen Yu can't change Lu's momentum without having her own momentum changed as well.

A cannon

example 2

- ▷ A cannon of mass 1000 kg fires a 10-kg shell at a velocity of 200 m/s. At what speed does the cannon recoil?
- ▷ The law of conservation of momentum tells us that

$$p_{\text{cannon},i} + p_{\text{shell},i} = p_{\text{cannon},f} + p_{\text{shell},f}.$$

Choosing a coordinate system in which the cannon points in the positive direction, the given information is

$$\begin{aligned} p_{\text{cannon},i} &= 0 \\ p_{\text{shell},i} &= 0 \\ p_{\text{shell},f} &= 2000 \text{ kg} \cdot \text{m/s}. \end{aligned}$$

We must have $p_{\text{cannon},f} = -2000 \text{ kg} \cdot \text{m/s}$, so the recoil velocity of the cannon is 2 m/s.

Ion drive

example 3

- ▷ The experimental solar-powered ion drive of the Deep Space 1 space probe expels its xenon gas exhaust at a speed of 30,000 m/s, ten times faster than the exhaust velocity for a typical chemical-fuel rocket engine. Roughly how many times greater is the maximum speed this spacecraft can reach, compared with a chemical-fueled probe with the same mass of fuel ("reaction mass") available for pushing out the back as exhaust?
- ▷ Momentum equals mass multiplied by velocity. Both spacecraft are assumed to have the same amount of reaction mass, and the ion drive's exhaust has a velocity ten times greater, so the momentum of its exhaust is ten times greater. Before the engine starts firing, neither the probe nor the exhaust has any momentum, so the total momentum of the system is zero. By conservation of momentum, the total momentum must also be zero after

all the exhaust has been expelled. If we define the positive direction as the direction the spacecraft is going, then the negative momentum of the exhaust is canceled by the positive momentum of the spacecraft. The ion drive allows a final speed that is ten times greater. (This simplified analysis ignores the fact that the reaction mass expelled later in the burn is not moving backward as fast, because of the forward speed of the already-moving spacecraft.)

3.1.2 Nonmechanical momentum

So far, it sounds as though conservation of momentum can be proved mathematically, unlike conservation of mass and energy, which are entirely based on observations. The proof, however, was only for a mechanical system, with interactions of the form $U(r)$. Conservation of momentum can be extended to other systems as well, but this generalization is based on experiments, not mathematical proof. Light is the most important example of momentum that doesn't equal mv — light doesn't have mass at all, but it does have momentum. For example, a flashlight left on for an hour would absorb about 10^{-5} kg · m/s of momentum as it recoiled.

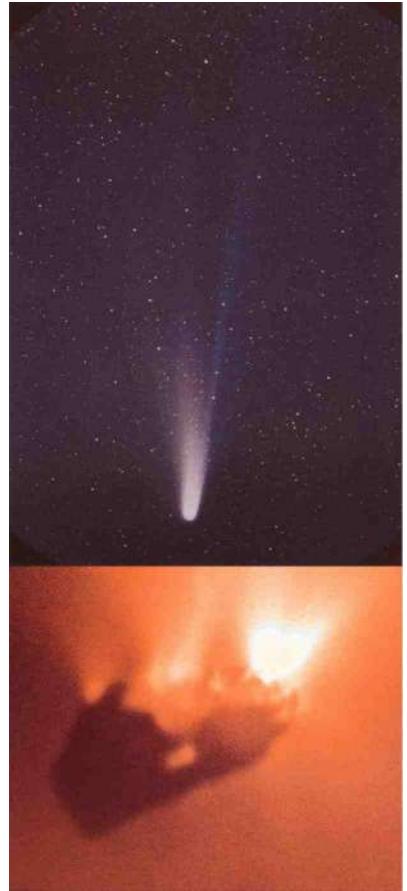
Halley's comet

example 4

Momentum is not always equal to mv . Halley's comet, shown in figure d, has a very elongated elliptical orbit, like those of many other comets. About once per century, its orbit brings it close to the sun. The comet's head, or nucleus, is composed of dirty ice, so the energy deposited by the intense sunlight gradually removes ice from the surface and turns it into water vapor. The bottom photo shows a view of the water coming off of the nucleus from the European Giotto space probe, which passed within 596 km of the comet's head on March 13, 1986.

The sunlight does not just carry energy, however. It also carries momentum. Once the steam comes off, the momentum of the sunlight impacting on it pushes it away from the sun, forming a tail as shown in the top image. The tail always points away from the sun, so when the comet is receding from the sun, the tail is in front. By analogy with matter, for which momentum equals mv , you would expect that massless light would have zero momentum, but the equation $p = mv$ is not the correct one for light, and light does have momentum. (Some comets also have a second tail, which is propelled by electrical forces rather than by the momentum of sunlight.)

The reason for bringing this up is not so that you can plug numbers into formulas in these exotic situations. The point is that the conservation laws have proven so sturdy exactly because they can easily be amended to fit new circumstances. The momentum of light will be a natural consequence of the discussion of the theory of relativity in chapter 7.



d / Halley's comet. Top: A photograph made from earth. Bottom: A view of the nucleus from the Giotto space probe. (W. Liller and European Space Agency)

3.1.3 Momentum compared to kinetic energy

Momentum and kinetic energy are both measures of the quantity of motion, and a sideshow in the Newton-Leibniz controversy over who invented calculus was an argument over whether mv (i.e., momentum) or mv^2 (i.e., kinetic energy without the $1/2$ in front) was the “true” measure of motion. The modern student can certainly be excused for wondering why we need both quantities, when their complementary nature was not evident to the greatest minds of the 1700s. The following table highlights their differences.

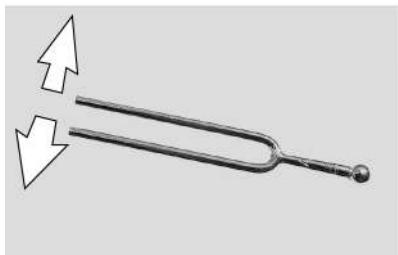
Kinetic energy...	Momentum...
doesn't depend on direction.	depends on direction.
is always positive, and cannot cancel out.	cancels with momentum in the opposite direction.
can be traded for forms of energy that do not involve motion. Kinetic energy is not a conserved quantity by itself.	is always conserved in a closed system.
is quadrupled if the velocity is doubled.	is doubled if the velocity is doubled.

Here are some examples that show the different behaviors of the two quantities.

A spinning top

example 5

A spinning top has zero total momentum, because for every moving point, there is another point on the opposite side that cancels its momentum. It does, however, have kinetic energy.



e / Examples 5 and 6. The momenta cancel, but the energies don't.

Why a tuning fork has two prongs

example 6

A tuning fork is made with two prongs so that they can vibrate in opposite directions, canceling their momenta. In a hypothetical version with only one prong, the momentum would have to oscillate, and this momentum would have to come from somewhere, such as the hand holding the fork. The result would be that vibrations would be transmitted to the hand and rapidly die out. In a two-prong fork, the two momenta cancel, but the energies don't.

Momentum and kinetic energy in firing a rifle

example 7

The rifle and bullet have zero momentum and zero kinetic energy to start with. When the trigger is pulled, the bullet gains some momentum in the forward direction, but this is canceled by the rifle's backward momentum, so the total momentum is still zero. The kinetic energies of the gun and bullet are both positive numbers, however, and do not cancel. The total kinetic energy is allowed to

increase, because kinetic energy is being traded for other forms of energy. Initially there is chemical energy in the gunpowder. This chemical energy is converted into heat, sound, and kinetic energy. The gun's "backward" kinetic energy does not refrigerate the shooter's shoulder!

The wobbly earth

example 8

As the moon completes half a circle around the earth, its motion reverses direction. This does not involve any change in kinetic energy. The reversed velocity does, however, imply a reversed momentum, so conservation of momentum in the closed earth-moon system tells us that the earth must also reverse its momentum. In fact, the earth wobbles in a little "orbit" about a point below its surface on the line connecting it and the moon. The two bodies' momenta always point in opposite directions and cancel each other out.

The earth and moon get a divorce

example 9

Why can't the moon suddenly decide to fly off one way and the earth the other way? It is not forbidden by conservation of momentum, because the moon's newly acquired momentum in one direction could be canceled out by the change in the momentum of the earth, supposing the earth headed off in the opposite direction at the appropriate, slower speed. The catastrophe is forbidden by conservation of energy, because their energies would have to increase greatly.

Momentum and kinetic energy of a glacier

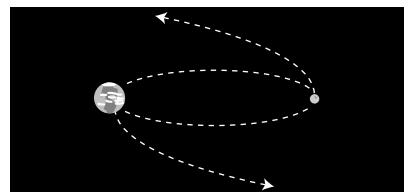
example 10

A cubic-kilometer glacier would have a mass of about 10^{12} kg. If it moves at a speed of 10^{-5} m/s, then its momentum is 10^7 kg · m/s. This is the kind of heroic-scale result we expect, perhaps the equivalent of the space shuttle taking off, or all the cars in LA driving in the same direction at freeway speed. Its kinetic energy, however, is only 50 J, the equivalent of the calories contained in a poppy seed or the energy in a drop of gasoline too small to be seen without a microscope. The surprisingly small kinetic energy is because kinetic energy is proportional to the square of the velocity, and the square of a small number is an even smaller number.

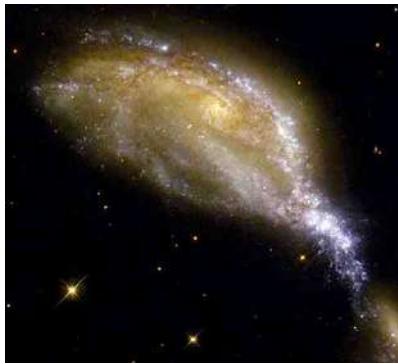
Discussion Questions

A If all the air molecules in the room settled down in a thin film on the floor, would that violate conservation of momentum? Conservation of energy?

B A refrigerator has coils in the back that get hot, and heat is molecular motion. These moving molecules have both energy and momentum. Why doesn't the refrigerator need to be tied to the wall to keep it from recoiling from the momentum it loses out the back?



f / Example 9.



g / This Hubble Space Telescope photo shows a small galaxy (yellow blob in the lower right) that has collided with a larger galaxy (spiral near the center), producing a wave of star formation (blue track) due to the shock waves passing through the galaxies' clouds of gas. This is considered a collision in the physics sense, even though it is statistically certain that no star in either galaxy ever struck a star in the other — the stars are very small compared to the distances between them. (NASA)

3.1.4 Collisions in one dimension

Physicists employ the term “collision” in a broader sense than in ordinary usage, applying it to any situation where objects interact for a certain period of time. A bat hitting a baseball, a cosmic ray damaging DNA, and a gun and a bullet going their separate ways are all examples of collisions in this sense. Physical contact is not even required. A comet swinging past the sun on a hyperbolic orbit is considered to undergo a collision, even though it never touches the sun. All that matters is that the comet and the sun interacted gravitationally with each other.

The reason for broadening the term “collision” in this way is that all of these situations can be attacked mathematically using the same conservation laws in similar ways. In our first example, conservation of momentum is all that is required.

Getting rear-ended

example 11

- ▷ Ms. Chang is rear-ended at a stop light by Mr. Nelson, and sues to make him pay her medical bills. He testifies that he was only going 55 km per hour when he hit Ms. Chang. She thinks he was going much faster than that. The cars skidded together after the impact, and measurements of the length of the skid marks show that their joint velocity immediately after the impact was 30 km per hour. Mr. Nelson's Nissan has a mass of 1400 kg, and Ms. Chang's Cadillac is 2400 kg. Is Mr. Nelson telling the truth?
▷ Since the cars skidded together, we can write down the equation for conservation of momentum using only two velocities, v for Mr. Nelson's velocity before the crash, and v' for their joint velocity afterward:

$$m_N v = m_N v' + m_C v'.$$

Solving for the unknown, v , we find

$$\begin{aligned}v &= \left(1 + \frac{m_C}{m_N}\right) v' \\&= 80 \text{ km/hr.}\end{aligned}$$

He is lying.

The above example was simple because both cars had the same velocity afterward. In many one-dimensional collisions, however, the two objects do not stick. If we wish to predict the result of such a collision, conservation of momentum does not suffice, because both velocities after the collision are unknown, so we have one equation in two unknowns.

Conservation of energy can provide a second equation, but its application is not as straightforward, because kinetic energy is only the particular form of energy that has to do with motion. In many collisions, part of the kinetic energy that was present before the collision is used to create heat or sound, or to break the objects

or permanently bend them. Cars, in fact, are carefully designed to crumple in a collision. Crumpling the car uses up energy, and that's good because the goal is to get rid of all that kinetic energy in a relatively safe and controlled way. At the opposite extreme, a superball is "super" because it emerges from a collision with almost all its original kinetic energy, having only stored it briefly as interatomic electrical energy while it was being squashed by the impact.

Collisions of the superball type, in which almost no kinetic energy is converted to other forms of energy, can thus be analyzed more thoroughly, because they have $K_f = K_i$, as opposed to the less useful inequality $K_f < K_i$ for a case like a tennis ball bouncing on grass. These two types of collisions are referred to, respectively, as elastic and inelastic. The extreme inelastic case is discussed further on p. 148.

Pool balls colliding head-on *example 12*

- ▷ Two pool balls collide head-on, so that the collision is restricted to one dimension. Pool balls are constructed so as to lose as little kinetic energy as possible in a collision, so under the assumption that no kinetic energy is converted to any other form of energy, what can we predict about the results of such a collision?
- ▷ Pool balls have identical masses, so we use the same symbol m for both. Conservation of energy and no loss of kinetic energy give us the two equations

$$\begin{aligned} mv_{1i} + mv_{2i} &= mv_{1f} + mv_{2f} \\ \frac{1}{2}mv_{1i}^2 + \frac{1}{2}mv_{2i}^2 &= \frac{1}{2}mv_{1f}^2 + \frac{1}{2}mv_{2f}^2 \end{aligned}$$

The masses and the factors of $1/2$ can be divided out, and we eliminate the cumbersome subscripts by replacing the symbols v_{1i}, \dots with the symbols A, B, C , and D :

$$\begin{aligned} A + B &= C + D \\ A^2 + B^2 &= C^2 + D^2. \end{aligned}$$

A little experimentation with numbers shows that given values of A and B , it is impossible to find C and D that satisfy these equations unless C and D equal A and B , or C and D are the same as A and B but swapped around. A formal proof of this fact is given in the sidebar. In the special case where ball 2 is initially at rest, this tells us that ball 1 is stopped dead by the collision, and ball 2 heads off at the velocity originally possessed by ball 1. This behavior will be familiar to players of pool.

Often, as in example 12, the details of the algebra are the least interesting part of the problem, and considerable physical insight can be gained simply by counting the number of unknowns and comparing to the number of equations. Suppose a beginner at pool

Gory details of the proof in example 12

The equation $A + B = C + D$ says that the change in one ball's velocity is equal and opposite to the change in the other's. We invent a symbol $x = C - A$ for the change in ball 1's velocity. The second equation can then be rewritten as $A^2 + B^2 = (A + x)^2 + (B - x)^2$. Squaring out the quantities in parentheses and then simplifying, we get $0 = Ax - Bx + x^2$. The equation has the trivial solution $x = 0$, i.e., neither ball's velocity is changed, but this is physically impossible because the balls can't travel through each other like ghosts. Assuming $x \neq 0$, we can divide by x and solve for $x = B - A$. This means that ball 1 has gained an amount of velocity exactly sufficient to match ball 2's initial velocity, and vice-versa. The balls must have swapped velocities.

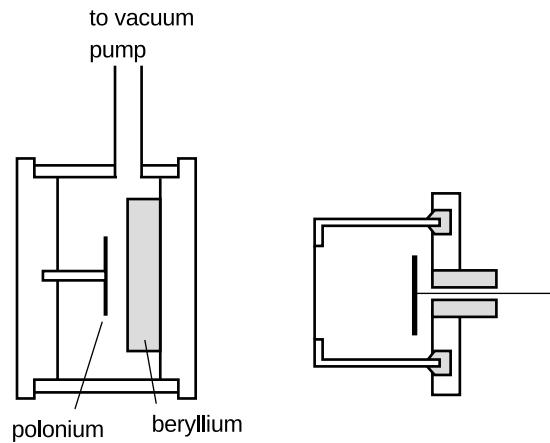
notices a case where her cue ball hits an initially stationary ball and stops dead. “Wow, what a good trick,” she thinks. “I bet I could never do that again in a million years.” But she tries again, and finds that she can’t help doing it even if she doesn’t want to. Luckily she has just learned about collisions in her physics course. Once she has written down the equations for conservation of momentum and no loss of kinetic energy, she really doesn’t have to complete the algebra. She knows that she has two equations in two unknowns, so there must be a well-defined solution. Once she has seen the result of one such collision, she knows that the same thing must happen every time. The same thing would happen with colliding marbles or croquet balls. It doesn’t matter if the masses or velocities are different, because that just multiplies both equations by some constant factor.

The discovery of the neutron

This was the type of reasoning employed by James Chadwick in his 1932 discovery of the neutron. At the time, the atom was imagined to be made out of two types of fundamental particles, protons and electrons. The protons were far more massive, and clustered together in the atom’s core, or nucleus. Electrical attraction caused the electrons to orbit the nucleus in circles, in much the same way that gravity kept the planets from cruising out of the solar system. Experiments showed, for example, that twice as much energy was required to strip the last electron off of a helium atom as was needed to remove the single electron from a hydrogen atom, and this was explained by saying that helium had two protons to hydrogen’s one. The trouble was that according to this model, helium would have two electrons and two protons, giving it precisely twice the mass of a hydrogen atom with one of each. In fact, helium has about four times the mass of hydrogen.

Chadwick suspected that the helium nucleus possessed two additional particles of a new type, which did not participate in electrical interactions at all, i.e., were electrically neutral. If these particles had very nearly the same mass as protons, then the four-to-one mass ratio of helium and hydrogen could be explained. In 1930, a new type of radiation was discovered that seemed to fit this description. It was electrically neutral, and seemed to be coming from the nuclei of light elements that had been exposed to other types of radiation. At this time, however, reports of new types of particles were a dime a dozen, and most of them turned out to be either clusters made of previously known particles or else previously known particles with higher energies. Many physicists believed that the “new” particle that had attracted Chadwick’s interest was really a previously known particle called a gamma ray, which was electrically neutral. Since gamma rays have no mass, Chadwick decided to try to determine the new particle’s mass and see if it was nonzero and

approximately equal to the mass of a proton.



h / Chadwick's subatomic pool table. A disk of the naturally occurring metal polonium provides a source of radiation capable of kicking neutrons out of the beryllium nuclei. The type of radiation emitted by the polonium is easily absorbed by a few mm of air, so the air has to be pumped out of the left-hand chamber. The neutrons, Chadwick's mystery particles, penetrate matter far more readily, and fly out through the wall and into the chamber on the right, which is filled with nitrogen or hydrogen gas. When a neutron collides with a nitrogen or hydrogen nucleus, it kicks it out of its atom at high speed, and this recoiling nucleus then rips apart thousands of other atoms of the gas. The result is an electrical pulse that can be detected in the wire on the right. Physicists had already calibrated this type of apparatus so that they could translate the strength of the electrical pulse into the velocity of the recoiling nucleus. The whole apparatus shown in the figure would fit in the palm of your hand, in dramatic contrast to today's giant particle accelerators.

Unfortunately a subatomic particle is not something you can just put on a scale and weigh. Chadwick came up with an ingenious solution. The masses of the nuclei of the various chemical elements were already known, and techniques had already been developed for measuring the speed of a rapidly moving nucleus. He therefore set out to bombard samples of selected elements with the mysterious new particles. When a direct, head-on collision occurred between a mystery particle and the nucleus of one of the target atoms, the nucleus would be knocked out of the atom, and he would measure its velocity.

Suppose, for instance, that we bombard a sample of hydrogen atoms with the mystery particles. Since the participants in the collision are fundamental particles, there is no way for kinetic energy to be converted into heat or any other form of energy, and Chadwick thus had two equations in three unknowns:

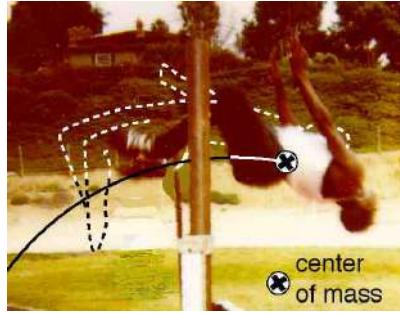
equation #1: conservation of momentum

equation #2: no loss of kinetic energy

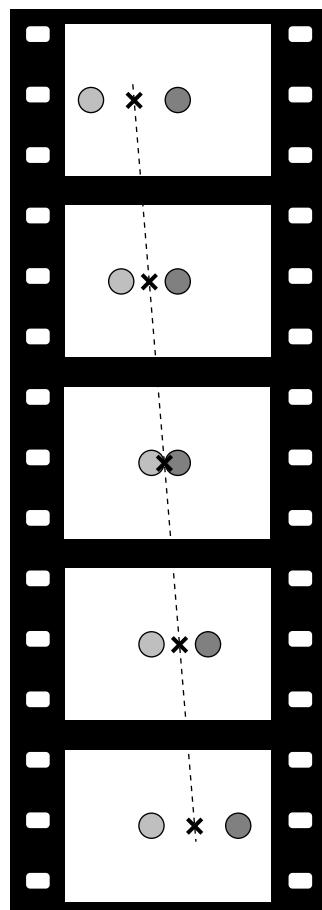
unknown #1: mass of the mystery particle

unknown #2: initial velocity of the mystery particle

unknown #3: final velocity of the mystery particle



i / The highjumper's body passes over the bar, but his center of mass passes under it. (*Dunia Young*)



j / Two pool balls collide.

The number of unknowns is greater than the number of equations, so there is no unique solution. But by creating collisions with nuclei of another element, nitrogen, he gained two more equations at the expense of only one more unknown:

equation #3: conservation of momentum in the new collision

equation #4: no loss of kinetic energy in the new collision

unknown #4: final velocity of the mystery particle in the new collision

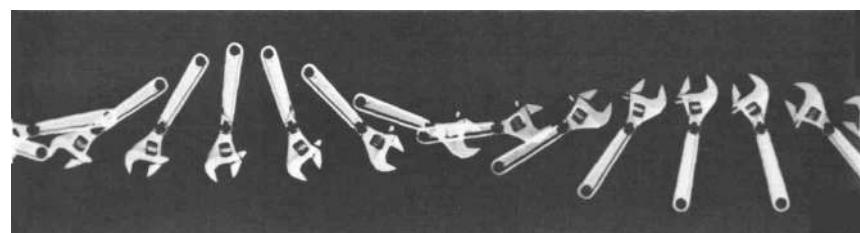
He was thus able to solve for all the unknowns, including the mass of the mystery particle, which was indeed within 1% of the mass of a proton. He named the new particle the neutron, since it is electrically neutral.

Discussion Questions

A Good pool players learn to make the cue ball spin, which can cause it not to stop dead in a head-on collision with a stationary ball. If this does not violate the laws of physics, what hidden assumption was there in the example in the text where it was proved that the cue ball must stop?

3.1.5 The center of mass

Figures i and k show two examples where a motion that appears complicated actually has a very simple feature. In both cases, there is a particular point, called the center of mass, whose motion is surprisingly simple. The highjumper flexes his body as he passes over the bar, so his motion is intrinsically very complicated, and yet his center of mass's motion is a simple parabola, just like the parabolic arc of a pointlike particle. The wrench's center of mass travels in a straight line as seen from above, which is what we'd expect for a pointlike particle flying through the air.



k / In this multiple-flash photograph, we see the wrench from above as it flies through the air, rotating as it goes. Its center of mass, marked with the black cross, travels along a straight line, unlike the other points on the wrench, which execute loops. (*PSSC Physics*)

The highjumper and the wrench are both complicated systems, each consisting of zillions of subatomic particles. To understand what's going on, let's instead look at a nice simple system, two pool balls colliding. We assume the balls are a closed system (i.e., their interaction with the felt surface is not important) and that their rotation is unimportant, so that we'll be able to treat each one as a single particle. By symmetry, the only place their center of mass can be is half-way in between, at an x coordinate equal to the average of the two balls' positions, $x_{cm} = (x_1 + x_2)/2$.

Figure j makes it appear that the center of mass, marked with an \times , moves with constant velocity to the right, regardless of the collision, and we can easily prove this using conservation of momentum:

$$\begin{aligned} v_{cm} &= dx_{cm}/dt \\ &= \frac{1}{2}(v_1 + v_2) \\ &= \frac{1}{2m}(mv_1 + mv_2) \\ &= \frac{p_{total}}{m_{total}} \end{aligned}$$

Since momentum is conserved, the last expression is constant, which proves that v_{cm} is constant.

Rearranging this a little, we have $p_{total} = m_{total}v_{cm}$. In other words, the total momentum of the system is the same as if all its mass was concentrated at the center of mass point.

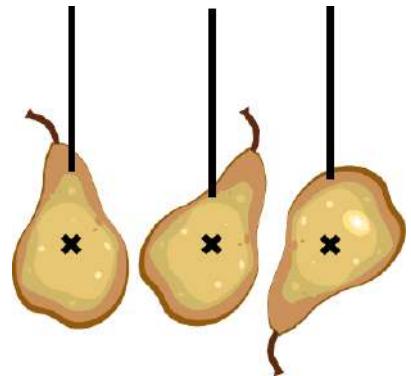
Sigma notation

When there is a large, potentially unknown number of particles, we can write sums like the ones occurring above using symbols like “+ . . .,” but that gets awkward. It's more convenient to use the Greek uppercase sigma, Σ , to indicate addition. For example, the sum $1^2 + 2^2 + 3^2 + 4^2 = 30$ could be written as

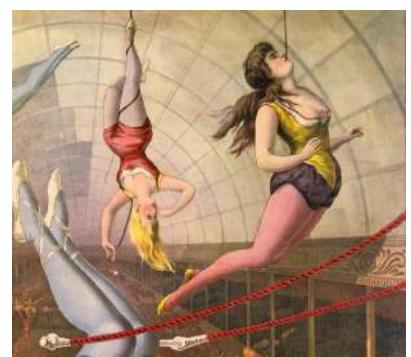
$$\sum_{j=1}^n j^2 = 30,$$

read “the sum from $j = 1$ to n of j^2 .” The variable j is a dummy variable, just like the dx in an integral that tells you you're integrating with respect to x , but has no significance outside the integral. The j below the sigma tells you what variable is changing from one term of the sum to the next, but j has no significance outside the sum.

As an example, let's generalize the proof of $p_{total} = m_{total}v_{cm}$ to the case of an arbitrary number n of identical particles moving in one dimension, rather than just two particles. The center of mass



l / No matter what point you hang the pear from, the string lines up with the pear's center of mass. The center of mass can therefore be defined as the intersection of all the lines made by hanging the pear in this way. Note that the X in the figure should not be interpreted as implying that the center of mass is on the surface — it is actually inside the pear.



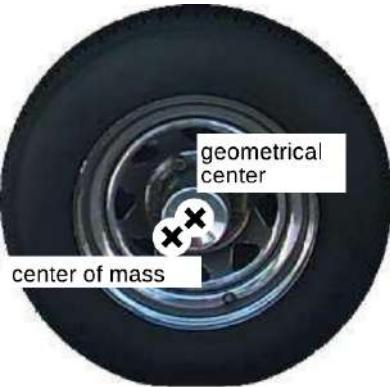
m / The circus performers hang with the ropes passing through their centers of mass.

is at

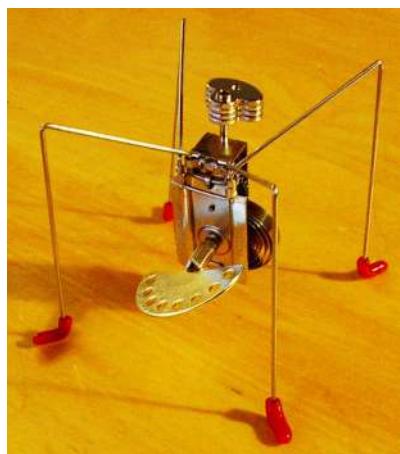
$$x_{cm} = \frac{1}{n} \sum_{j=1}^n x_j,$$

where x_1 is the mass of the first particle, and so on. The velocity of the center of mass is

$$\begin{aligned} v_{cm} &= dx_{cm}/dt \\ &= \frac{1}{n} \sum_{j=1}^n v_j \\ &= \frac{1}{nm} \sum_{j=1}^n mv_j \\ &= \frac{p_{total}}{m_{total}} \end{aligned}$$



n / An improperly balanced wheel has a center of mass that is not at its geometric center. When you get a new tire, the mechanic clamps little weights to the rim to balance the wheel.



o / This toy was intentionally designed so that the mushroom-shaped piece of metal on top would throw off the center of mass. When you wind it up, the mushroom spins, but the center of mass doesn't want to move, so the rest of the toy tends to counter the mushroom's motion, causing the whole thing to jump around.

What about a system containing objects with unequal masses, or containing more than two objects? The reasoning above can be generalized to a weighted average:

$$x_{cm} = \frac{\sum_{j=1}^n m_j x_j}{\sum_{j=1}^n m_j}$$

The solar system's center of mass

example 13

In the discussion of the sun's gravitational field on page 99, I mentioned in a footnote that the sun doesn't really stay in one place while the planets orbit around it. Actually, motion is relative, so it's meaningless to ask whether the sun is absolutely at rest, but it is meaningful to ask whether it moves in a straight line at constant velocity. We can now see that since the solar system is a closed system, its total momentum must be constant, and $p_{total} = m_{total}v_{cm}$ then tells us that it's the solar system's center of mass that has constant velocity, not the sun. The sun wobbles around this point irregularly due to its interactions with the planets, Jupiter in particular.

The earth-moon system

example 14

The earth-moon system is much simpler than the solar system because it contains only two objects. Where is the center of mass of this system? Let $x=0$ be the earth's center, so that the moon lies at $x = 3.8 \times 10^5$ km. Then

$$\begin{aligned} x_{cm} &= \frac{\sum_{j=1}^2 m_j x_j}{\sum_{j=1}^2 m_j} \\ &= \frac{m_1 x_1 + m_2 x_2}{m_1 + m_2}, \end{aligned}$$

and letting 1 be the earth and 2 the moon, we have

$$x_{cm} = \frac{m_{earth} \times 0 + m_{moon}x_{moon}}{m_{earth} + m_{moon}}$$
$$= 4600 \text{ km},$$

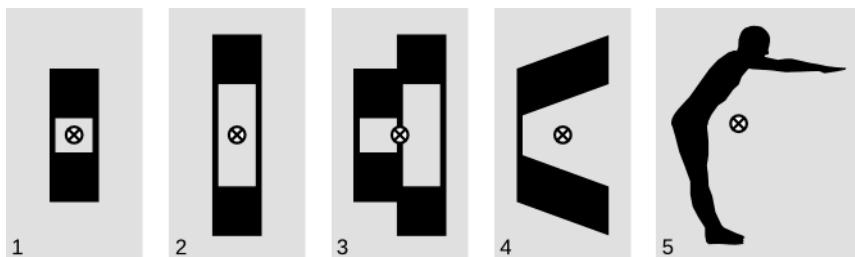
or about three quarters of the way from the earth's center to its surface.

The center of mass as an average

example 15

- ▷ Explain how we know that the center of mass of each object is at the location shown in figure p.

p / Example 15.



▷ The center of mass is a sort of average, so the height of the centers of mass in 1 and 2 has to be midway between the two squares, because that height is the average of the heights of the two squares. Example 3 is a combination of examples 1 and 2, so we can find its center of mass by averaging the horizontal positions of their centers of mass. In example 4, each square has been skewed a little, but just as much mass has been moved up as down, so the average vertical position of the mass hasn't changed. Example 5 is clearly not all that different from example 4, the main difference being a slight clockwise rotation, so just as in example 4, the center of mass must be hanging in empty space, where there isn't actually any mass. Horizontally, the center of mass must be between the heels and toes, or else it wouldn't be possible to stand without tipping over.

Momentum and Galilean relativity

example 16

The principle of Galilean relativity states that the laws of physics are supposed to be equally valid in all inertial frames of reference. If we first calculate some momenta in one frame of reference and find that momentum is conserved, and then rework the whole problem in some other frame of reference that is moving with respect to the first, the numerical values of the momenta will all be different. Even so, momentum will still be conserved. All that matters is that we work a single problem in one consistent frame of reference.

One way of proving this is to apply the equation $p_{total} = m_{total}v_{cm}$. If the velocity of one frame relative to the other is u , then the only effect of changing frames of reference is to change v_{cm} from its original value to $v_{cm} + u$. This adds a constant onto the momentum, which has no effect on conservation of momentum.

self-check A

The figure shows a gymnast holding onto the inside of a big wheel. From inside the wheel, how could he make it roll one way or the other?

- ▷ Answer, p. 1059



q / Self-check A.

3.1.6 The center of mass frame of reference

A particularly useful frame of reference in many cases is the frame that moves along with the center of mass, called the center of mass (c.m.) frame. In this frame, the total momentum is zero. The following examples show how the center of mass frame can be a powerful tool for simplifying our understanding of collisions.

A collision of pool balls viewed in the c.m. frame example 17

If you move your head so that your eye is always above the point halfway in between the two pool balls, as in figure r, you are viewing things in the center of mass frame. In this frame, the balls come toward the center of mass at equal speeds. By symmetry, they must therefore recoil at equal speeds along the lines on which they entered. Since the balls have essentially swapped paths in the center of mass frame, the same must also be true in any other frame. This is the same result that required laborious algebra to prove previously without the concept of the center of mass frame.

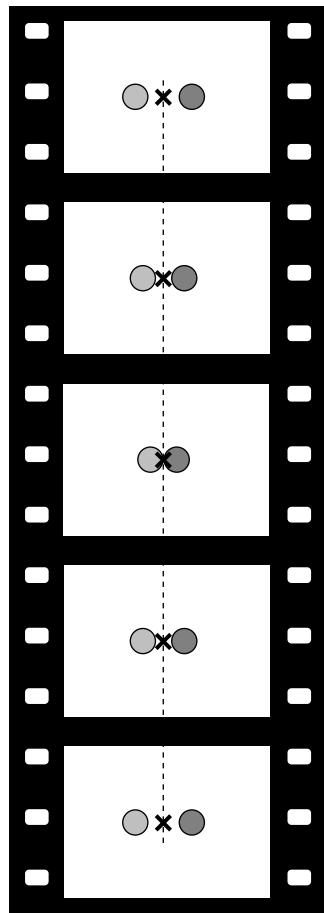
The slingshot effect example 18

It is a counterintuitive fact that a spacecraft can pick up speed by swinging around a planet, if it arrives in the opposite direction compared to the planet's motion. Although there is no physical contact, we treat the encounter as a one-dimensional collision, and analyze it in the center of mass frame. Since Jupiter is so much more massive than the spacecraft, the center of mass is essentially fixed at Jupiter's center, and Jupiter has zero velocity in the center of mass frame, as shown in figure 3.1.6. The c.m. frame is moving to the left compared to the sun-fixed frame used in figure 3.1.6, so the spacecraft's initial velocity is greater in this frame than in the sun's frame.

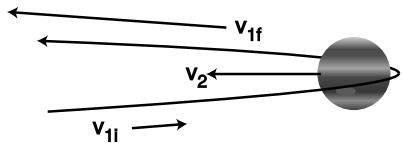
Things are simpler in the center of mass frame, because it is more symmetric. In the sun-fixed frame, the incoming leg of the encounter is rapid, because the two bodies are rushing toward each other, while their separation on the outbound leg is more gradual, because Jupiter is trying to catch up. In the c.m. frame, Jupiter is sitting still, and there is perfect symmetry between the incoming and outgoing legs, so by symmetry we have $v_{1f} = -v_{1i}$. Going back to the sun-fixed frame, the spacecraft's final velocity is increased by the frames' motion relative to each other. In the sun-fixed frame, the spacecraft's velocity has increased greatly.

Einstein's motorcycle example 19

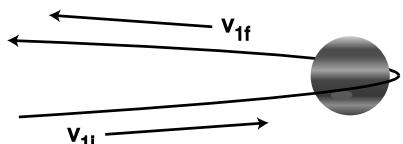
We've assumed we were dealing with a system of material objects, for which the equation $p = mv$ was true. What if our system contains only light rays, or a mixture of light and matter? As a college student, Einstein kept worrying about what a beam of light would look like if you could ride alongside it on a motorcycle. In other words, he imagined putting himself in the light beam's



r / The same collision of two pool balls, but now seen in the center of mass frame of reference.



s / The sun's frame of reference.



t / The c.m. frame.

center of mass frame. Chapter 7 discusses Einstein's resolution of this problem, but the basic point is that you *can't* ride the motorcycle alongside the light beam, because material objects can't go as fast as the speed of light. A beam of light has no center of mass frame of reference.

Discussion Questions

A Make up a numerical example of two unequal masses moving in one dimension at constant velocity, and verify the equation $p_{total} = m_{total}v_{cm}$ over a time interval of one second.

B A more massive tennis racquet or baseball bat makes the ball fly off faster. Explain why this is true, using the center of mass frame. For simplicity, assume that the racquet or bat is simply sitting still before the collision, and that the hitter's hands do not make any force large enough to have a significant effect over the short duration of the impact.

3.1.7 Totally inelastic collisions

On p. 139 we discussed collisions that were totally elastic (no conversion of KE into other types of energy). A useful application of the center of mass frame of reference is to the description of the opposite extreme, a totally *inelastic* collision.

A totally inelastic collision cannot just be defined as one in which all the KE is converted into other forms, both because the definition would depend on our frame of reference and because there is a constraint imposed by conservation of momentum. Let's say that a golfer hits a ball. In the frame of reference of the grass, it would violate conservation of momentum if the ball were to stay put while the club simply stopped moving. If such a complete cessation of motion is to happen, then it must occur in the center of mass frame of reference. In the c.m. frame, there is zero total momentum both before and after the collision. Thus if we observe no motion at all after the collision, we must be in the c.m. frame.

Therefore we define a totally inelastic collision as one in which there is no motion in the c.m. frame in the final state. An observer watching such a collision, in any frame, will see that the amount of KE transformed into other forms of energy is as great as possible subject to conservation of momentum.

When objects touch physically (possibly crumpling or changing shape during the collision) in a totally elastic collision, the final state in the c.m. frame is one in which the two objects are at rest and touching. In other frames of reference, we see the objects stick to each other and travel away together after the collision. An example of this type was example 11 on p. 138, in which one car rear-ended another, and they stuck together as a unit after the crash.

3.2 Force in one dimension

3.2.1 Momentum transfer

For every conserved quantity, we can define an associated rate of flow. An open system can have mass transferred in or out of it, and we can measure the rate of mass flow, dm/dt in units of kg/s. Energy can flow in or out, and the rate of energy transfer is the power, $P = dE/dt$, measured in watts.³ The rate of *momentum* transfer is called force,

$$F = \frac{dp}{dt} \quad [\text{definition of force}].$$

The units of force are $\text{kg}\cdot\text{m}/\text{s}^2$, which can be abbreviated as newtons, $1 \text{ N} = \text{kg}\cdot\text{m}/\text{s}^2$. Newtons are unfortunately not as familiar as watts. A newton is about how much force you'd use to pet a dog. The most powerful rocket engine ever built, the first stage of the Saturn V that sent astronauts to the moon, had a thrust of about 30 million newtons. In one dimension, positive and negative signs indicate the direction of the force — a positive force is one that pushes or pulls in the direction of the positive x axis.

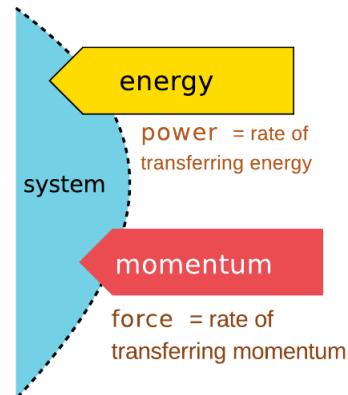
Walking into a lamppost

example 20

- ▷ Starting from rest, you begin walking, bringing your momentum up to $100 \text{ kg}\cdot\text{m}/\text{s}$. You walk straight into a lamppost. Why is the momentum change of $-100 \text{ kg}\cdot\text{m}/\text{s}$ so much more painful than the change of $+100 \text{ kg}\cdot\text{m}/\text{s}$ when you started walking?
- ▷ The forces are not really constant, but for this type of qualitative discussion we can pretend they are, and approximate dp/dt as $\Delta p/\Delta t$. It probably takes you about 1 s to speed up initially, so the ground's force on you is $F = \Delta p/\Delta t \approx 100 \text{ N}$. Your impact with the lamppost, however, is over in the blink of an eye, say $1/10 \text{ s}$ or less. Dividing by this much smaller Δt gives a much larger force, perhaps thousands of newtons (with a negative sign).

This is also the principle of airbags in cars. The time required for the airbag to decelerate your head is fairly long: the time it takes your face to travel 20 or 30 cm. Without an airbag, your face would have hit the dashboard, and the time interval would have been the much shorter time taken by your skull to move a couple of centimeters while your face compressed. Note that either way, the same amount of momentum is transferred: the entire momentum of your head.

Force is defined as a derivative, and the derivative of a sum is the sum of the derivatives. Therefore force is additive: when more than one force acts on an object, you add the forces to find out what happens. An important special case is that forces can cancel. Consider your body sitting in a chair as you read this book. Let



a / Power and force are the rates at which energy and momentum are transferred.



b / The airbag increases Δt so as to reduce $F = \Delta p/\Delta t$.

³Recall that uppercase P is power, while lowercase p is momentum.

the positive x axis be upward. The chair's upward force on you is represented with a positive number, which cancels out with the earth's downward gravitational force, which is negative. The total rate of momentum transfer into your body is zero, and your body doesn't change its momentum.

Finding momentum from force

example 21

- ▷ An object of mass m starts at rest at $t = t_0$. A force varying as $F = bt^{-2}$, where b is a constant, begins acting on it. Find the greatest speed it will ever have.

▷

$$\begin{aligned} F &= \frac{dp}{dt} \\ dp &= F dt \\ p &= \int F dt + p_0 \\ &= -\frac{b}{t} + p_0, \end{aligned}$$

where p_0 is a constant of integration. The given initial condition is that $p = 0$ at $t = t_0$, so we find that $p_0 = b/t_0$. The negative term gets closer to zero with increasing time, so the maximum momentum is achieved by letting t approach infinity. That is, the object will never stop speeding up, but it will also never surpass a certain speed. In the limit $t \rightarrow \infty$, we identify p_0 as the momentum that the object will approach asymptotically. The maximum velocity is $v = p_0/m = b/mt_0$.

Discussion Question

A Many collisions, like the collision of a bat with a baseball, appear to be instantaneous. Most people also would not imagine the bat and ball as bending or being compressed during the collision. Consider the following possibilities:

- (1) The collision is instantaneous.
- (2) The collision takes a finite amount of time, during which the ball and bat retain their shapes and remain in contact.
- (3) The collision takes a finite amount of time, during which the ball and bat are bending or being compressed.

How can two of these be ruled out based on energy or momentum considerations?

3.2.2 Newton's laws

Although momentum is the third conserved quantity we've encountered, historically it was the first to be discovered. Isaac Newton formulated a complete treatment of mechanical systems in terms of force and momentum. Newton's theory was based on three laws of motion, which we now think of as consequences of conservation of mass, energy, and momentum.



c / Isaac Newton (1643-1727).

Newton's laws in one dimension:

Newton's first law: If there is no force acting on an object, it stays in the same state of motion.

Newton's second law: The sum of all the forces acting on an object determines the rate at which its momentum changes, $F_{total} = dp/dt$.

Newton's third law: Forces occur in opposite pairs. If object A interacts with object B, then A's force on B and B's force on A are related by $F_{AB} = -F_{BA}$.

The second law is the definition of force, which we've already encountered.⁴ The first law is a special case of the second law — if dp/dt is zero, then $p = mv$ is a constant, and since mass is conserved, constant p implies constant v . The third law is a restatement of conservation of momentum: for two objects interacting, we have constant total momentum, so $0 = \frac{d}{dt}(p_A + p_B) = F_{BA} + F_{AB}$.

a=F/m

example 22

Many modern textbooks restate Newton's second law as $a = F/m$, i.e., as an equation that predicts an object's acceleration based on the force exerted on it. This is easily derived from Newton's original form as follows: $a = dv/dt = (dp/dt)/m = F/m$.

Gravitational force related to g

example 23

As a special case of the previous example, consider an object in free fall, and let the x axis point down. Then $a = +g$, and $F = ma = mg$. For example, the gravitational force on a 1 kg mass at the earth's surface is about 9.8 N. Even if other forces act on the object, and it isn't in free fall, the gravitational force on it is still the same, and can still be calculated as mg .

Changing frames of reference

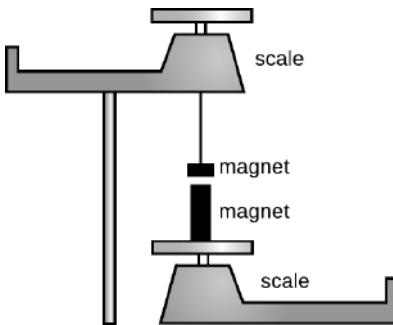
example 24

Suppose we change from one frame reference into another, which is moving relative to the first one at a constant velocity u . If an object of mass m is moving at velocity v (which need not be constant), then the effect is to change its momentum from mv in one frame to $mv+mu$ in the other. Force is defined as the derivative of momentum with respect to time, and the derivative of a constant is zero, so adding the constant mu has no effect on the result. We therefore conclude that observers in different inertial frames of reference agree on forces.

Using the third law correctly

If you've already accepted Galilean relativity in your heart, then there is nothing really difficult about the first and second laws. The third law, however, is more of a conceptual challenge. The first

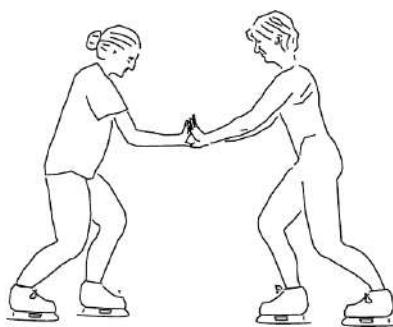
⁴This is with the benefit of hindsight. At the time, the word “force” already had certain connotations, and people thought they understood what it meant and how to measure it, e.g., by using a spring scale. From their point of view, $F = dp/dt$ was not a definition but a testable — and controversial! — statement.



d / Two magnets exert forces on each other.



e / It doesn't make sense for the man to talk about the woman's money canceling out his bar tab, because there is no good reason to combine his debts and her assets.



f / Newton's third law does not mean that forces always cancel out so that nothing can ever move. If these two ice skaters, initially at rest, push against each other, they will both move.

hurdle is that it is counterintuitive. Is it really true that if a fighter jet collides with a mosquito, the mosquito's force on the jet is just as strong as the jet's force on the mosquito? Yes, it is true, but it is hard to believe at first. That amount of force simply has more of an effect on the mosquito, because it has less mass.

A more humane and practical experiment is shown in figure d. A large magnet and a small magnet are weighed separately, and then one magnet is hung from the pan of the top balance so that it is directly above the other magnet. There is an attraction between the two magnets, causing the reading on the top scale to increase and the reading on the bottom scale to decrease. The large magnet is more "powerful" in the sense that it can pick up a heavier paperclip from the same distance, so many people have a strong expectation that one scale's reading will change by a far different amount than the other. Instead, we find that the two changes are equal in magnitude but opposite in direction, so the upward force of the top magnet on the bottom magnet is of the same magnitude as the downward force of the bottom magnet on the top magnet.

To students, it often sounds as though Newton's third law implies nothing could ever change its motion, since the two equal and opposite forces would always cancel. As illustrated in figure e, the fallacy arises from assuming that we can add things that it doesn't make sense to add. It only makes sense to add up forces that are acting on the same object, whereas two forces related to each other by Newton's third law are always acting on two different objects. If two objects are interacting via a force and no other forces are involved, then *both* objects will accelerate — in opposite directions, as shown in figure f!

Here are some suggestions for avoiding misapplication of Newton's third law:

1. It always relates exactly two forces, not more.
2. The two forces involve exactly two objects, in the pattern A on B, B on A.
3. The two forces are always of the same type, e.g., friction and friction, or gravity and gravity.

Directions of forces

We've already seen that momentum, unlike energy, has a direction in space. Since force is defined in terms of momentum, force also has a direction in space. For motion in one dimension, we have to pick a coordinate system, and given that choice, forces and momenta will be positive or negative. We've already used signs to represent directions of forces in Newton's third law, $F_{AB} = -F_{BA}$.

There is, however, a complication with force that we were able to avoid with momentum. If an object is moving on a line, we're guaranteed that its momentum is in one of two directions: the two directions along the line. But even an object that stays on a line may still be subject to forces that act perpendicularly to the line. For example, suppose a coin is sliding to the right across a table, \mathbf{h} , and let's choose a positive x axis that points to the right. The coin's motion is along a horizontal line, and its momentum is positive and decreasing. Because the momentum is decreasing, its time derivative $d\mathbf{p}/dt$ is negative. This derivative equals the horizontal force of friction F_1 , and its negative sign tells us that this force on the coin is to the left.

But there are also vertical forces on the coin. The Earth exerts a downward gravitational force F_2 on it, and the table makes an upward force F_3 that prevents the coin from sinking into the wood. In fact, without these vertical forces the horizontal frictional force wouldn't exist: surfaces don't exert friction against one another unless they are being pressed together.

To avoid mathematical complication, we want to postpone the full three-dimensional treatment of force and momentum until section 3.4. For now, we'll limit ourselves to examples like the coin, in which the motion is confined to a line, and any forces perpendicular to the line cancel each other out.

Discussion Questions

A Criticize the following incorrect statement:

"If an object is at rest and the total force on it is zero, it stays at rest. There can also be cases where an object is moving and keeps on moving without having any total force on it, but that can only happen when there's no friction, like in outer space."

B The table gives laser timing data for Ben Johnson's 100 m dash at the 1987 World Championship in Rome. (His world record was later revoked because he tested positive for steroids.) How does the total force on him change over the duration of the race?

C You hit a tennis ball against a wall. Explain any and all incorrect ideas in the following description of the physics involved: "According to Newton's third law, there has to be a force opposite to your force on the ball. The opposite force is the ball's mass, which resists acceleration, and also air resistance."

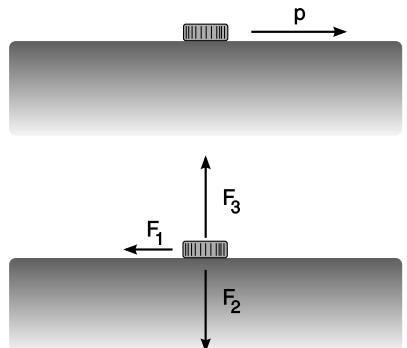
D Tam Anh grabs Sarah by the hand and tries to pull her. She tries to remain standing without moving. A student analyzes the situation as follows. "If Tam Anh's force on Sarah is greater than her force on him, he can get her to move. Otherwise, she'll be able to stay where she is." What's wrong with this analysis?

3.2.3 What force is not

Violin teachers have to endure their beginning students' screeching. A frown appears on the woodwind teacher's face as she watches



g / A swimmer doing the breast stroke pushes backward against the water. By Newton's third law, the water pushes forward on him.



h / A coin slides across a table. Even for motion in one dimension, some of the forces may not be along the line of the motion.

30	3.80
40	4.67
50	5.53
60	6.38
70	7.23
80	8.10
90	8.96
100	9.83

Discussion question B.

her student take a breath with an expansion of his ribcage but none in his belly. What makes physics teachers cringe is their students' verbal statements about forces. Below I have listed several dicta about what force is not.

Force is not a property of one object.

A great many of students' incorrect descriptions of forces could be cured by keeping in mind that a force is an interaction of two objects, not a property of one object.

Incorrect statement: "That magnet has a lot of force."

X If the magnet is one millimeter away from a steel ball bearing, they may exert a very strong attraction on each other, but if they were a meter apart, the force would be virtually undetectable. The magnet's strength can be rated using certain electrical units (ampere – meters²), but not in units of force.

Force is not a measure of an object's motion.

If force is not a property of a single object, then it cannot be used as a measure of the object's motion.

Incorrect statement: "The freight train rumbled down the tracks with awesome force."

X Force is not a measure of motion. If the freight train collides with a stalled cement truck, then some awesome forces will occur, but if it hits a fly the force will be small.

Force is not energy.

Incorrect statement: "How can my chair be making an upward force on my rear end? It has no power!"

X Power is a concept related to energy, e.g., a 100-watt lightbulb uses up 100 joules per second of energy. When you sit in a chair, no energy is used up, so forces can exist between you and the chair without any need for a source of power.

Force is not stored or used up.

Because energy can be stored and used up, people think force also can be stored or used up.

Incorrect statement: "If you don't fill up your tank with gas, you'll run out of force."

X Energy is what you'll run out of, not force.

Forces need not be exerted by living things or machines.

Transforming energy from one form into another usually requires some kind of living or mechanical mechanism. The concept is not applicable to forces, which are an interaction between objects, not a thing to be transferred or transformed.

Incorrect statement: "How can a wooden bench be making an upward force on my rear end? It doesn't have any springs or anything inside it."

X No springs or other internal mechanisms are required. If the bench didn't make any force on you, you would obey Newton's second law and fall through it. Evidently it does make a force on you!

A force is the direct cause of a change in motion.

I can click a remote control to make my garage door change from being at rest to being in motion. My finger's force on the button, however, was not the force that acted on the door. When we speak of a force on an object in physics, we are talking about a force that acts directly. Similarly, when you pull a reluctant dog along by its leash, the leash and the dog are making forces on each other, not your hand and the dog. The dog is not even touching your hand.

self-check B

Which of the following things can be correctly described in terms of force?

- (1) A nuclear submarine is charging ahead at full steam.
- (2) A nuclear submarine's propellers spin in the water.
- (3) A nuclear submarine needs to refuel its reactor periodically. ▷

Answer, p. 1059

Discussion Questions

A Criticize the following incorrect statement: "If you shove a book across a table, friction takes away more and more of its force, until finally it stops."

B You hit a tennis ball against a wall. Explain any and all incorrect ideas in the following description of the physics involved: "The ball gets some force from you when you hit it, and when it hits the wall, it loses part of that force, so it doesn't bounce back as fast. The muscles in your arm are the only things that a force can come from."

3.2.4 Forces between solids

Conservation laws are more fundamental than Newton's laws, and they apply where Newton's laws don't, e.g., to light and to the internal structure of atoms. However, there are certain problems that are much easier to solve using Newton's laws. As a trivial example, if you drop a rock, it could conserve momentum and energy by levitating, or by falling in the usual manner.⁵ With Newton's laws, however, we can reason that $a = F/m$, so the rock must respond to the gravitational force by accelerating.

Less trivially, suppose a person is hanging onto a rope, and we want to know if she will slip. Unlike the case of the levitating rock, here the no-motion solution could be perfectly reasonable if her grip is strong enough. We know that her hand's interaction with the rope is fundamentally an electrical interaction between the atoms in the surface of her palm and the nearby atoms in the surface of the rope.

⁵This pathological solution was first noted on page 83, and discussed in more detail on page 1025.

For practical problem-solving, however, this is a case where we're better off forgetting the fundamental classification of interactions at the atomic level and working with a more practical, everyday classification of forces. In this practical scheme, we have three types of forces that can occur between solid objects in contact:

<i>A normal force, F_n,</i>	is perpendicular to the surface of contact, and prevents objects from passing through each other by becoming as strong as necessary (up to the point where the objects break). “Normal” means perpendicular.
<i>Static friction, F_s,</i>	is parallel to the surface of contact, and prevents the surfaces from starting to slip by becoming as strong as necessary, up to a maximum value of $F_{s,max}$. “Static” means not moving, i.e., not slipping.
<i>Kinetic friction, F_k,</i>	is parallel to the surface of contact, and tends to slow down any slippage once it starts. “Kinetic” means moving, i.e., slipping.

self-check C

Can a frictionless surface exert a normal force? Can a frictional force exist without a normal force?

▷ Answer, p. 1059

If you put a coin on this page, which is horizontal, gravity pulls down on the coin, but the atoms in the paper and the coin repel each other electrically, and the atoms are compressed until the repulsion becomes strong enough to stop the downward motion of the coin. We describe this complicated and invisible atomic process by saying that the paper makes an upward normal force on the coin, and the coin makes a downward normal force on the paper. (The two normal forces are related by Newton's third law. In fact, Newton's third law only relates forces that are of the same type.)

If you now tilt the book a little, static friction keeps the coin from slipping. The picture at the microscopic level is even more complicated than the previous description of the normal force. One model is to think of the tiny bumps and depressions in the coin as settling into the similar irregularities in the paper. This model predicts that rougher surfaces should have more friction, which is sometimes true but not always. Two very smooth, clean glass surfaces or very well finished machined metal surfaces may actually stick *better* than rougher surfaces would, the probable explanation being that there is some kind of chemical bonding going on, and the smoother surfaces allow more atoms to be in contact.

Finally, as you tilt the book more and more, there comes a point where static friction reaches its maximum value. The surfaces be-

come unstuck, and the coin begins to slide over the paper. Kinetic friction slows down this slipping motion significantly. In terms of energy, kinetic friction is converting mechanical energy into heat, just like when you rub your hands together to keep warm. One model of kinetic friction is that the tiny irregularities in the two surfaces bump against each other, causing vibrations whose energy rapidly converts to heat and sound — you can hear this sound if you rub your fingers together near your ear.

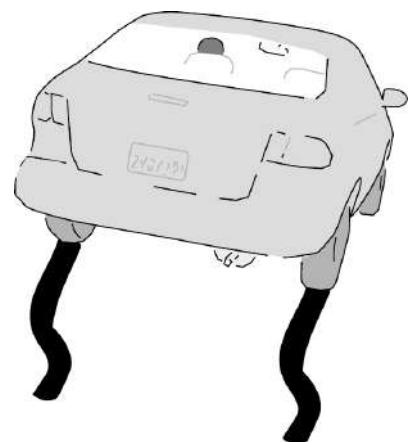
For *dry* surfaces, experiments show that the following equations usually work fairly well:



i / Static friction: the tray doesn't slip on the waiter's fingers.

$$F_{s,max} \approx \mu_s F_n,$$

and



j / Kinetic friction: the car skids.

$$F_k \approx \mu_k F_n,$$

where μ_s , the coefficient of static friction, and μ_k , the coefficient of kinetic friction, are constants that depend on the properties of the two surfaces, such as what they're made of and how rough they are.

self-check D

1. When a baseball player slides in to a base, is the friction static, or kinetic?
2. A mattress stays on the roof of a slowly accelerating car. Is the friction static, or kinetic?
3. Does static friction create heat? Kinetic friction? ▷ Answer, p. 1059

Maximum acceleration of a car

example 25

▷ Rubber on asphalt gives $\mu_k \approx 0.4$ and $\mu_s \approx 0.6$. What is the upper limit on a car's acceleration on a flat road, assuming that the engine has plenty of power and that air friction is negligible?

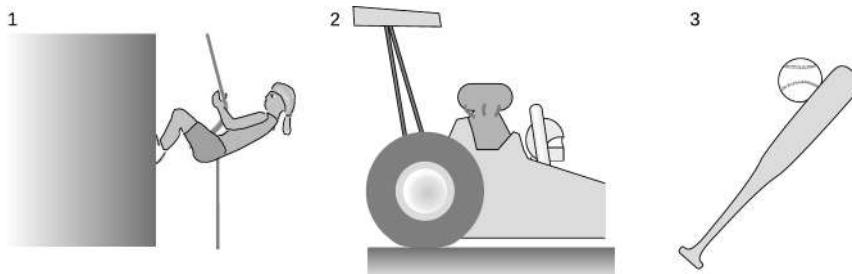
▷ This isn't a flying car, so we don't expect it to accelerate vertically. The vertical forces acting on the car should cancel out. The earth makes a downward gravitational force on the car whose absolute value is mg , so the road apparently makes an upward normal force of the same magnitude, $F_n = mg$.

Now what about the horizontal motion? As is always true, the coefficient of static friction is greater than the coefficient of kinetic friction, so the maximum acceleration is obtained with static friction, i.e., the driver should try not to burn rubber. The maximum force of static friction is $F_{s,max} = \mu_s F_n = \mu_s mg$. The maximum acceleration is $a = F_s/m = \mu_s g \approx 6 \text{ m/s}^2$. This is true regardless of how big the tires are, since the experimentally determined relationship $F_{s,max} = \mu_s F_n$ is independent of surface area.

self-check E

Find the direction of each of the forces in figure k. ▷ Answer, p. 1059

k / 1. The cliff's normal force on the climber's feet. 2. The track's static frictional force on the wheel of the accelerating dragster. 3. The ball's normal force on the bat.



Locomotives

example 26

Looking at a picture of a locomotive, l, we notice two obvious things that are different from an automobile. Where a car typically has two drive wheels, a locomotive normally has many — ten in this example. (Some also have smaller, unpowered wheels in front of and behind the drive wheels, but this example doesn't.) Also, cars these days are generally built to be as light as possible for their size, whereas locomotives are very massive, and no effort seems to be made to keep their weight low. (The steam locomotive in the photo is from about 1900, but this is true even for modern diesel and electric trains.)

The reason locomotives are built to be so heavy is for traction. The upward normal force of the rails on the wheels, F_N , cancels the downward force of gravity, F_W , so ignoring plus and minus signs, these two forces are equal in absolute value, $F_N = F_W$.



Given this amount of normal force, the maximum force of static friction is $F_s = \mu_s F_N = \mu_s F_W$. This static frictional force, of the rails pushing forward on the wheels, is the only force that can accelerate the train, pull it uphill, or cancel out the force of air resistance while cruising at constant speed. The coefficient of static friction for steel on steel is about 1/4, so no locomotive can pull with a force greater than about 1/4 of its own weight. If the engine is capable of supplying more than that amount of force, the result will be simply to break static friction and spin the wheels.

The reason this is all so different from the situation with a car is that a car isn't pulling something else. If you put extra weight in a car, you improve the traction, but you also increase the inertia of the car, and make it just as hard to accelerate. In a train, the inertia is almost all in the cars being pulled, not in the locomotive.

The other fact we have to explain is the large number of driving wheels. First, we have to realize that increasing the number of driving wheels neither increases nor decreases the total amount of static friction, because static friction is independent of the amount of surface area in contact. (The reason four-wheel-drive is good in a car is that if one or more of the wheels is slipping on ice or in mud, the other wheels may still have traction. This isn't typically an issue for a train, since all the wheels experience the same conditions.) The advantage of having more driving wheels on a train is that it allows us to increase the weight of the locomotive without crushing the rails, or damaging bridges.

3.2.5 Fluid friction

Try to drive a nail into a waterfall and you will be confronted with the main difference between solid friction and fluid friction. Fluid friction is purely kinetic; there is no static fluid friction. The nail in the waterfall may tend to get dragged along by the water flowing past it, but it does not stick in the water. The same is true for gases such as air: recall that we are using the word "fluid" to include both gases and liquids.

Unlike kinetic friction between solids, fluid friction increases rapidly with velocity. It also depends on the shape of the object, which is why a fighter jet is more streamlined than a Model T. For objects of the same shape but different sizes, fluid friction typically scales up with the cross-sectional area of the object, which is one of the main reasons that an SUV gets worse mileage on the freeway.



m / The wheelbases of the Hummer H3 and the Toyota Prius are surprisingly similar, differing by only 10%. The main difference in shape is that the Hummer is much taller and wider. It presents a much greater cross-sectional area to the wind, and this is the main reason that it uses about 2.5 times more gas on the freeway.

than a compact car.

Discussion Question

A Criticize the following analysis: “A book is sitting on a table. I shove it, overcoming static friction. Then it slows down until it has less force than static friction, and it stops.”

3.2.6 Analysis of forces

Newton’s first and second laws deal with the total of all the forces exerted on a specific object, so it is very important to be able to figure out what forces there are. Once you have focused your attention on one object and listed the forces on it, it is also helpful to describe all the corresponding forces that must exist according to Newton’s third law. We refer to this as “analyzing the forces” in which the object participates.

A barge

example 27

A barge is being pulled along a canal by teams of horses on the shores. Analyze all the forces in which the barge participates.

<i>force acting on barge</i>	<i>force related to it by Newton’s third law</i>
ropes’ forward normal forces on barge	barge’s backward normal force on ropes
water’s backward fluid friction force on barge	barge’s forward fluid friction force on water
planet earth’s downward gravitational force on barge	barge’s upward gravitational force on earth
water’s upward “floating” force on barge	barge’s downward “floating” force on water

Here I’ve used the word “floating” force as an example of a sensible invented term for a type of force not classified on the tree in the previous section. A more formal technical term would be “hydrostatic force.”

Note how the pairs of forces are all structured as “A’s force on B, B’s force on A”: ropes on barge and barge on ropes; water on barge and barge on water. Because all the forces in the left column are forces acting on the barge, all the forces in the right column are forces being exerted by the barge, which is why each entry in the column begins with “barge.”

Often you may be unsure whether you have missed one of the forces. Here are three strategies for checking your list:

1. See what physical result would come from the forces you’ve found so far. Suppose, for instance, that you’d forgotten the “floating” force on the barge in the example above. Looking at the forces you’d found, you would have found that there was a downward gravitational force on the barge which was not canceled by any upward force. The barge isn’t supposed to sink, so you know you need to find a fourth, upward force.
2. Whenever one solid object touches another, there will be a normal force, and possibly also a frictional force; check for both.
3. Make a drawing of the object, and draw a dashed boundary line around it that separates it from its environment. Look for points on the boundary where other objects come in contact with your object. This strategy guarantees that you’ll find every contact force that acts on the object, although it won’t

help you to find non-contact forces.

The following is another example in which we can profit by checking against our physical intuition for what should be happening.

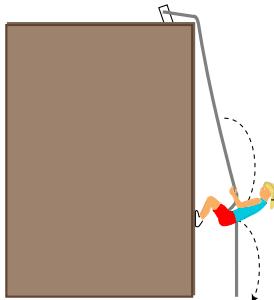
Rappelling

example 28

As shown in the figure below, Cindy is rappelling down a cliff. Her downward motion is at constant speed, and she takes little hops off of the cliff, as shown by the dashed line. Analyze the forces in which she participates at a moment when her feet are on the cliff and she is pushing off.

<i>force acting on Cindy</i>	<i>force related to it by Newton's third law</i>
planet earth's downward gravitational force on Cindy	Cindy's upward gravitational force on earth
ropes upward frictional force on Cindy (her hand)	Cindy's downward frictional force on the rope
cliff's rightward normal force on Cindy	Cindy's leftward normal force on the cliff

The two vertical forces cancel, which is what they should be doing if she is to go down at a constant rate. The only horizontal force on her is the cliff's force, which is not canceled by any other force, and which therefore will produce an acceleration of Cindy to the right. This makes sense, since she is hopping off. (This solution is a little oversimplified, because the rope is slanting, so it also applies a small leftward force to Cindy. As she flies out to the right, the slant of the rope will increase, pulling her back in more strongly.)



I believe that constructing the type of table described in this section is the best method for beginning students. Most textbooks, however, prescribe a pictorial way of showing all the forces acting on an object. Such a picture is called a free-body diagram. It should not be a big problem if a future physics professor expects you to be able to draw such diagrams, because the conceptual reasoning is the same. You simply draw a picture of the object, with arrows representing the forces that are acting on it. Arrows representing contact forces are drawn from the point of contact, noncontact forces from the center of mass. Free-body diagrams do not show the equal and opposite forces exerted by the object itself.

Discussion Questions

A When you fire a gun, the exploding gases push outward in all directions, causing the bullet to accelerate down the barrel. What Newton's-third-law pairs are involved? [Hint: Remember that the gases themselves are an object.]

B In the example of the barge going down the canal, I referred to a "floating" or "hydrostatic" force that keeps the boat from sinking. If you were adding a new branch on the force-classification tree to represent this force, where would it go?

C A pool ball is rebounding from the side of the pool table. Analyze the forces in which the ball participates during the short time when it is in contact with the side of the table.

D The earth's gravitational force on you, i.e., your weight, is always equal to mg , where m is your mass. So why can you get a shovel to go deeper into the ground by jumping onto it? Just because you're jumping, that doesn't mean your mass or weight is any greater, does it?

3.2.7 Transmission of forces by low-mass objects

You're walking your dog. The dog wants to go faster than you do, and the leash is taut. Does Newton's third law guarantee that your force on your end of the leash is equal and opposite to the dog's force on its end? If they're not exactly equal, is there any reason why they should be approximately equal?

If there was no leash between you, and you were in direct contact with the dog, then Newton's third law would apply, but Newton's third law cannot relate your force on the leash to the dog's force on the leash, because that would involve three separate objects. Newton's third law only says that your force on the leash is equal and opposite to the leash's force on you,

$$F_{yL} = -F_{Ly},$$

and that the dog's force on the leash is equal and opposite to its force on the dog

$$F_{dL} = -F_{Ld}.$$

Still, we have a strong intuitive expectation that whatever force we make on our end of the leash is transmitted to the dog, and vice-versa. We can analyze the situation by concentrating on the forces that act on the leash, F_{dL} and F_{yL} . According to Newton's second law, these relate to the leash's mass and acceleration:

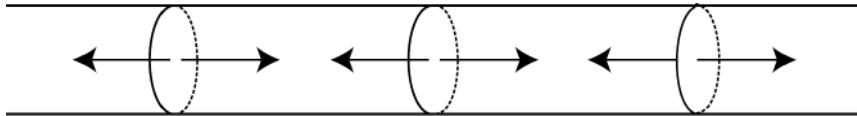
$$F_{dL} + F_{yL} = m_L a_L.$$

The leash is far less massive than any of the other objects involved, and if m_L is very small, then apparently the total force on the leash is also very small, $F_{dL} + F_{yL} \approx 0$, and therefore

$$F_{dL} \approx -F_{yL}.$$

Thus even though Newton's third law does not apply directly to these two forces, we can approximate the low-mass leash as if it was not intervening between you and the dog. It's at least approximately as if you and the dog were acting directly on each other, in which case Newton's third law would have applied.

In general, low-mass objects can be treated approximately as if they simply transmitted forces from one object to another. This can be true for strings, ropes, and cords, and also for rigid objects such as rods and sticks.



n / If we imagine dividing a taut rope up into small segments, then any segment has forces pulling outward on it at each end. If the rope is of negligible mass, then all the forces equal $+T$ or $-T$, where T , the tension, is a single number.

If you look at a piece of string under a magnifying glass as you pull on the ends more and more strongly, you will see the fibers straightening and becoming taut. Different parts of the string are apparently exerting forces on each other. For instance, if we think of the two halves of the string as two objects, then each half is exerting a force on the other half. If we imagine the string as consisting of many small parts, then each segment is transmitting a force to the next segment, and if the string has very little mass, then all the forces are equal in magnitude. We refer to the magnitude of the forces as the tension in the string, T .

The term “tension” refers only to internal forces within the string. If the string makes forces on objects at its ends, then those forces are typically normal or frictional forces (example 29).

Types of force made by ropes

example 29

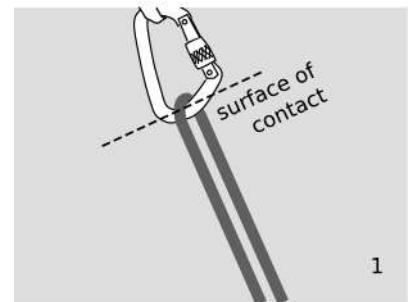
- ▷ Analyze the forces in figures o/1 and o/2.
- ▷ In all cases, a rope can only make “pulling” forces, i.e., forces that are parallel to its own length and that are toward itself, not away from itself. You can’t push with a rope!

In o/1, the rope passes through a type of hook, called a carabiner, used in rock climbing and mountaineering. Since the rope can only pull along its own length, the direction of its force on the carabiner must be down and to the right. This is perpendicular to the surface of contact, so the force is a normal force.

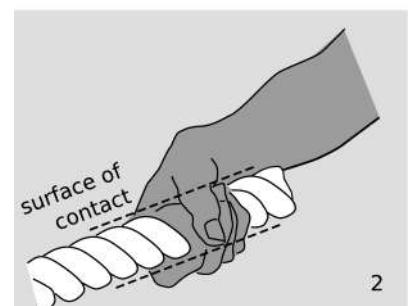
<i>force acting on carabiner</i>	<i>force related to it by Newton’s third law</i>
rope’s normal force on carabiner ↓	carabiner’s normal force on rope ↑

(There are presumably other forces acting on the carabiner from other hardware above it.)

In figure o/2, the rope can only exert a net force at its end that is parallel to itself and in the pulling direction, so its force on the hand is down and to the left. This is parallel to the surface of contact, so it must be a frictional force. If the rope isn’t slipping through the hand, we have static friction. Friction can’t exist with-



1



2

o / Example 29. The forces between the rope and other objects are normal and frictional forces.

out normal forces. These forces are perpendicular to the surface of contact. For simplicity, we show only two pairs of these normal forces, as if the hand were a pair of pliers.

<i>force acting on person</i>	<i>force related to it by Newton's third law</i>
rope's static frictional force on person ↖	person's static frictional force on rope ↗
rope's normal force on person ↖	person's normal force on rope ↘
rope's normal force on person ↘	person's normal force on rope ↖

(There are presumably other forces acting on the person as well, such as gravity.)

If a rope goes over a pulley or around some other object, then the tension throughout the rope is approximately equal so long as the pulley has negligible mass and there is not too much friction. A rod or stick can be treated in much the same way as a string, but it is possible to have either compression or tension.

Discussion Question

A When you step on the gas pedal, is your foot's force being transmitted in the sense of the word used in this section?

3.2.8 Work

Energy transferred to a particle

To change the kinetic energy, $K = (1/2)mv^2$, of a particle moving in one dimension, we must change its velocity. That will entail a change in its momentum, $p = mv$, as well, and since force is the rate of transfer of momentum, we conclude that the only way to change a particle's kinetic energy is to apply a force.⁶ A force in the same direction as the motion speeds it up, increasing the kinetic energy, while a force in the opposite direction slows it down.

Consider an infinitesimal time interval during which the particle moves an infinitesimal distance dx , and its kinetic energy changes by dK . In one dimension, we represent the direction of the force and the direction of the motion with positive and negative signs for F and dx , so the relationship among the signs can be summarized as follows:

⁶The converse isn't true, because kinetic energy doesn't depend on the direction of motion, but momentum does. We can change a particle's momentum without changing its energy, as when a pool ball bounces off a bumper, reversing the sign of p .

$F > 0$	$dx > 0$	$dK > 0$
$F < 0$	$dx < 0$	$dK > 0$
$F > 0$	$dx < 0$	$dK < 0$
$F < 0$	$dx > 0$	$dK < 0$

This looks exactly like the rule for determining the sign of a product, and we can easily show using the chain rule that this is indeed a multiplicative relationship:

$$\begin{aligned} dK &= \frac{dK}{dv} \frac{dv}{dt} \frac{dt}{dx} dx \quad [\text{chain rule}] \\ &= (mv)(a)(1/v) dx \\ &= m a dx \\ &= F dx \quad [\text{Newton's second law}] \end{aligned}$$

We can verify that force multiplied by distance has units of energy:

$$\begin{aligned} N \cdot m &= \frac{\text{kg} \cdot \text{m}/\text{s}}{\text{s}} \times \text{m} \\ &= \text{kg} \cdot \text{m}^2/\text{s}^2 \\ &= \text{J} \end{aligned}$$

A TV picture tube

example 30

- ▷ At the back of a typical TV's picture tube, electrical forces accelerate each electron to an energy of 5×10^{-16} J over a distance of about 1 cm. How much force is applied to a single electron? (Assume the force is constant.) What is the corresponding acceleration?

- ▷ Integrating

$$dK = F dx,$$

we find

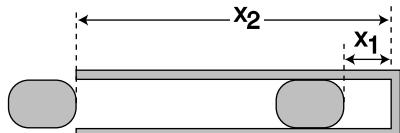
$$K_f - K_i = F(x_f - x_i)$$

or

$$\Delta K = F \Delta x.$$

The force is

$$\begin{aligned} F &= \Delta K / \Delta x \\ &= \frac{5 \times 10^{-16} \text{ J}}{0.01 \text{ m}} \\ &= 5 \times 10^{-14} \text{ N.} \end{aligned}$$



p / A simplified drawing of an airgun.

This may not sound like an impressive force, but it's enough to supply an electron with a spectacular acceleration. Looking up the mass of an electron on p. 1072, we find

$$a = F/m \\ = 5 \times 10^{16} \text{ m/s}^2.$$

An air gun

example 31

▷ An airgun, figure p, uses compressed air to accelerate a pellet. As the pellet moves from x_1 to x_2 , the air decompresses, so the force is not constant. Using methods from chapter 5, one can show that the air's force on the pellet is given by $F = bx^{-7/5}$. A typical high-end airgun used for competitive target shooting has

$$x_1 = 0.046 \text{ m}, \\ x_2 = 0.41 \text{ m},$$

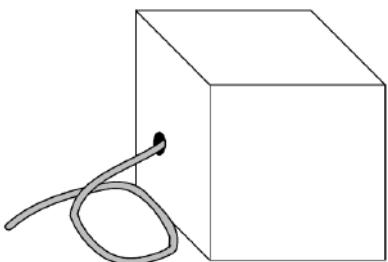
and

$$b = 4.4 \text{ N}\cdot\text{m}^{7/5}.$$

What is the kinetic energy of the pellet when it leaves the muzzle? (Assume friction is negligible.)

▷ Since the force isn't constant, it would be incorrect to do $F = \Delta K/\Delta x$. Integrating both sides of the equation $dK = F dx$, we have

$$\begin{aligned} \Delta K &= \int_{x_1}^{x_2} F dx \\ &= -\frac{5b}{2} \left(x_2^{-2/5} - x_1^{-2/5} \right) \\ &= 22 \text{ J} \end{aligned}$$



q / The black box does work by reeling in its cable.

In general, when energy is transferred by a force,⁷ we use the term *work* to refer to the amount of energy transferred. This is different from the way the word is used in ordinary speech. If you stand for a long time holding a bag of cement, you get tired, and everyone will agree that you've worked hard, but you haven't changed the energy of the cement, so according to the definition of the physics term, you haven't done any work on the bag. There has been an energy transformation inside your body, of chemical energy into heat, but this just means that one part of your body did positive work (lost energy) while another part did a corresponding amount of negative work (gained energy).

⁷The part of the definition about "by a force" is meant to exclude the transfer of energy by heat conduction, as when a stove heats soup.

Work in general

I derived the expression $F dx$ for one particular type of kinetic-energy transfer, the work done in accelerating a particle, and then defined work as a more general term. Is the equation correct for other types of work as well? For example, if a force lifts a mass m against the resistance of gravity at constant velocity, the increase in the mass's gravitational energy is $d(mgy) = mg dy = F dy$, so again the equation works, but this still doesn't prove that the equation is *always* correct as a way of calculating energy transfers.

Imagine a black box⁸, containing a gasoline-powered engine, which is designed to reel in a steel cable, exerting a certain force F . For simplicity, we imagine that this force is always constant, so we can talk about Δx rather than an infinitesimal dx . If this black box is used to accelerate a particle (or any mass without internal structure), and no other forces act on the particle, then the original derivation applies, and the work done by the box is $W = F\Delta x$. Since F is constant, the box will run out of gas after reeling in a certain amount of cable Δx . The chemical energy inside the box has decreased by $-W$, while the mass being accelerated has gained W worth of kinetic energy.⁹

Now what if we use the black box to pull a plow? The energy increase in the outside world is of a different type than before; it takes the forms of (1) the gravitational energy of the dirt that has been lifted out to the sides of the furrow, (2) frictional heating of the dirt and the plowshare, and (3) the energy needed to break up the dirt clods (a form of electrical energy involving the attractions among the atoms in the clod). The box, however, only communicates with the outside world via the hole through which its cable passes. The amount of chemical energy lost by the gasoline can therefore only depend on F and Δx , so it is the same $-W$ as when the box was being used to accelerate a mass, and thus by conservation of energy, the work done on the outside world is again W .

This is starting to sound like a proof that the force-times-distance method is always correct, but there was one subtle assumption involved, which was that the force was exerted at one point (the end of the cable, in the black box example). Real life often isn't like that. For example, a cyclist exerts forces on both pedals at once. Serious cyclists use toe-clips, and the conventional wisdom is that one should use equal amounts of force on the upstroke and downstroke, to make full use of both sets of muscles. This is a two-dimensional example, since the pedals go in circles. We're only discussing one-dimensional motion right now, so let's just pretend that the upstroke and down-

⁸"Black box" is a traditional engineering term for a device whose inner workings we don't care about.

⁹For conceptual simplicity, we ignore the transfer of heat energy to the outside world via the exhaust and radiator. In reality, the sum of these energies plus the useful kinetic energy transferred would equal W .



r / The wheel spinning in the air has $K_{cm} = 0$. The space shuttle has all its kinetic energy in the form of center of mass motion, $K = K_{cm}$. The rolling ball has some, but not all, of its energy in the form of center of mass motion, $K_{cm} < K$. (Space Shuttle photo by NASA)

stroke are both executed in straight lines. Since the forces are in opposite directions, one is positive and one is negative. The cyclist's *total* force on the crank set is zero, but the work done isn't zero. We have to add the work done by each stroke, $W = F_1\Delta x_1 + F_2\Delta x_2$. (I'm pretending that both forces are constant, so we don't have to do integrals.) Both terms are positive; one is a positive number multiplied by a positive number, while the other is a negative times a negative.

This might not seem like a big deal — just remember not to use the total force — but there are many situations where the total force is all we can measure. The ultimate example is heat conduction. Heat conduction is not supposed to be counted as a form of work, since it occurs without a force. But at the atomic level, there are forces, and work is done by one atom on another. When you hold a hot potato in your hand, the transfer of heat energy through your skin takes place with a total force that's extremely close to zero. At the atomic level, atoms in your skin are interacting electrically with atoms in the potato, but the attractions and repulsions add up to zero total force. It's just like the cyclist's feet acting on the pedals, but with zillions of forces involved instead of two. There is no practical way to measure all the individual forces, and therefore we can't calculate the total energy transferred.

To summarize, $\sum F_j dx_j$ is a correct way of calculating work, where F_j is the individual force acting on particle j , which moves a distance dx_j . However, this is only useful if you can identify all the individual forces and determine the distance moved at each point of contact. For convenience, I'll refer to this as the *work theorem*. (It doesn't have a standard name.)

There is, however, something useful we can do with the total force. We can use it to calculate the part of the work done on an object that consists of a change in the kinetic energy it has due to the motion of its center of mass. The proof is essentially the same as the proof on p. 165, except that now we don't assume the force is acting on a single particle, so we have to be a little more delicate. Let the object consist of n particles. Its total kinetic energy is $K = \sum_{j=1}^n (1/2)m_j v_j^2$, but this is what we've already realized *can't* be calculated using the total force. The kinetic energy it has due to motion of its center of mass is

$$K_{cm} = \frac{1}{2}m_{total}v_{cm}^2.$$

Figure r shows some examples of the distinction between K_{cm} and

K . Differentiating K_{cm} , we have

$$\begin{aligned} dK_{cm} &= m_{total}v_{cm} dv_{cm} \\ &= m_{total}v_{cm} \frac{dv_{cm}}{dt} \frac{dt}{dx_{cm}} dx_{cm} \quad [\text{chain rule}] \\ &= m_{total} \frac{dv_{cm}}{dt} dx_{cm} \quad [dt/dx_{cm} = 1/v_{cm}] \\ &= \frac{dp_{total}}{dt} dx_{cm} \quad [p_{total} = m_{total}v_{cm}] \\ &= F_{total} dx_{cm} \end{aligned}$$

I'll call this the *kinetic energy theorem* — like the work theorem, it has no standard name.

An ice skater pushing off from a wall

example 32

The kinetic energy theorem tells us how to calculate the skater's kinetic energy if we know the amount of force and the distance her center of mass travels while she is pushing off.

The work theorem tells us that the wall does no work on the skater, since the point of contact isn't moving. This makes sense, because the wall does not have any source of energy.

Absorbing an impact without recoiling?

example 33

▷ Is it possible to absorb an impact without recoiling? For instance, if a ping-pong ball hits a brick wall, does the wall "give" at all?

▷ There will always be a recoil. In the example proposed, the wall will surely have some energy transferred to it in the form of heat and vibration. The work theorem tells us that we can only have an energy transfer if the distance traveled by the point of contact is nonzero.

Dragging a refrigerator at constant velocity

example 34

The fridge's momentum is constant, so there is no net momentum transfer, and the total force on it must be zero: your force is canceling the floor's kinetic frictional force. The kinetic energy theorem is therefore true but useless. It tells us that there is zero total force on the refrigerator, and that the refrigerator's kinetic energy doesn't change.

The work theorem tells us that the work you do equals your hand's force on the refrigerator multiplied by the distance traveled. Since we know the floor has no source of energy, the only way for the floor and refrigerator to gain energy is from the work you do. We can thus calculate the total heat dissipated by friction in the refrigerator and the floor.

Note that there is no way to find how much of the heat is dissipated in the floor and how much in the refrigerator.

Accelerating a cart***example 35***

If you push on a cart and accelerate it, there are two forces acting on the cart: your hand's force, and the static frictional force of the ground pushing on the wheels in the opposite direction.

Applying the work theorem to your force tells us how to calculate the work you do.

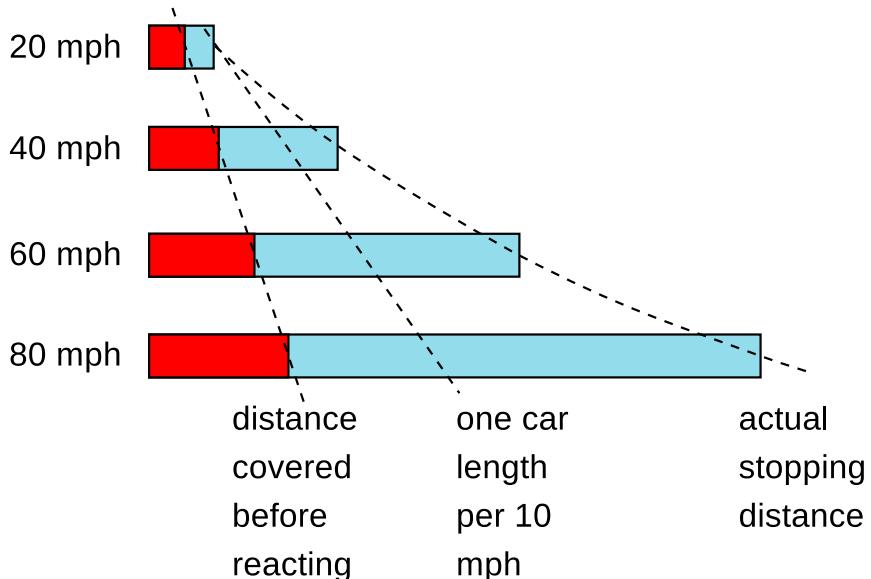
Applying the work theorem to the floor's force tells us that the floor does no work on the cart. There is no motion at the point of contact, because the atoms in the floor are not moving. (The atoms in the surface of the wheel are also momentarily at rest when they touch the floor.) This makes sense, because the floor does not have any source of energy.

The kinetic energy theorem refers to the total force, and because the floor's backward force cancels part of your force, the total force is less than your force. This tells us that only part of your work goes into the kinetic energy associated with the forward motion of the cart's center of mass. The rest goes into rotation of the wheels.

Discussion Questions

A Criticize the following incorrect statement: "A force doesn't do any work unless it's causing the object to move."

B To stop your car, you must first have time to react, and then it takes some time for the car to slow down. Both of these times contribute to the distance you will travel before you can stop. The figure shows how the average stopping distance increases with speed. Because the stopping distance increases more and more rapidly as you go faster, the rule of one car length per 10 m.p.h. of speed is not conservative enough at high speeds. In terms of work and kinetic energy, what is the reason for the more rapid increase at high speeds?



s / Discussion question B.

3.2.9 Simple Machines

Conservation of energy provided the necessary tools for analyzing some mechanical systems, such as the seesaw on page 85 and the pulley arrangements of the homework problems on page 122, but we could only analyze those machines by computing the total energy of the system. That approach wouldn't work for systems like the biceps/forearm machine on page 85, or the one in figure t, where the energy content of the person's body is impossible to compute directly. Even though the seesaw and the biceps/forearm system were clearly just two different forms of the lever, we had no way to treat them both on the same footing. We can now successfully attack such problems using the work and kinetic energy theorems.

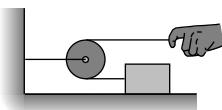
Constant tension around a pulley

example 36

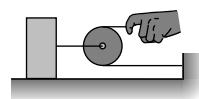
▷ In figure t, what is the relationship between the force applied by the person's hand and the force exerted on the block?

▷ If we assume the rope and the pulley are ideal, i.e., frictionless and massless, then there is no way for them to absorb or release energy, so the work done by the hand must be the same as the work done on the block. Since the hand and the block move the same distance, the work theorem tells us the two forces are the same.

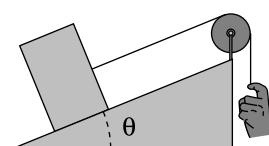
Similar arguments provide an alternative justification for the statement made in section 3.2.7 that show that an idealized rope exerts the same force, the tension, anywhere it's attached to something, and the same amount of force is also exerted by each segment of the rope on the neighboring segments. Going around an ideal pulley also has no effect on the tension.



t / The force is transmitted to the block.



u / A mechanical advantage of 2.



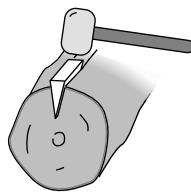
v / An inclined plane.

This is an example of a simple machine, which is any mechanical system that manipulates forces to do work. This particular machine reverses the direction of the motion, but doesn't change the force or the speed of motion.

A mechanical advantage

example 37

The idealized pulley in figure u has negligible mass, so its kinetic energy is zero, and the kinetic energy theorem tells us that the total force on it is zero. We know, as in the preceding example, that the two forces pulling it to the right are equal to each other, so the force on the left must be twice as strong. This simple machine doubles the applied force, and we refer to this ratio as a *mechanical advantage* (M.A.) of 2. There's no such thing as a free lunch, however; the distance traveled by the load is cut in half, and there is no increase in the amount of work done.



w / A wedge.

Inclined plane and wedge

example 38

In figure v, the force applied by the hand is equal to the one applied to the load, but there is a mechanical advantage compared to the force that would have been required to lift the load straight up. The distance traveled up the inclined plane is greater by a factor of $1/\sin \theta$, so by the work theorem, the force is smaller by a factor of $\sin \theta$, and we have $M.A.=1/\sin \theta$. The wedge, w, is similar.

Archimedes' screw

example 39

In one revolution, the crank travels a distance $2\pi b$, and the water rises by a height h . The mechanical advantage is $2\pi b/h$.

3.2.10 Force related to interaction energy

In section 2.3, we saw that there were two equivalent ways of looking at gravity, the gravitational field and the gravitational energy. They were related by the equation $dU = mg dr$, so if we knew the field, we could find the energy by integration, $U = \int mg dr$, and if we knew the energy, we could find the field by differentiation, $g = (1/m) dU / dr$.

The same approach can be applied to other interactions, for example a mass on a spring. The main difference is that only in gravitational interactions does the strength of the interaction depend on the mass of the object, so in general, it doesn't make sense to separate out the factor of m as in the equation $dU = mg dr$. Since $F = mg$ is the gravitational force, we can rewrite the equation in the more suggestive form $dU = F dr$. This form no longer refers to gravity specifically, and can be applied much more generally. The only remaining detail is that I've been fairly cavalier about positive and negative signs up until now. That wasn't such a big problem for gravitational interactions, since gravity is always attractive, but it requires more careful treatment for nongravitational forces, where we don't necessarily know the direction of the force in advance, and

x / Archimedes' screw

we need to use positive and negative signs carefully for the direction of the force.

In general, suppose that forces are acting on a particle — we can think of them as coming from other objects that are “off stage” — and that the interaction between the particle and the off-stage objects can be characterized by an interaction energy, U , which depends only on the particle’s position, x . Using the kinetic energy theorem, we have $dK = F dx$. (It’s not necessary to write K_{cm} , since a particle can’t have any other kind of kinetic energy.) Conservation of energy tells us $dK + dU = 0$, so the relationship between force and interaction energy is $dU = -F dx$, or

$$F = -\frac{dU}{dx} \quad [\text{relationship between force and interaction energy}].$$

Force exerted by a spring

example 40

▷ A mass is attached to the end of a spring, and the energy of the spring is $U = (1/2)kx^2$, where x is the position of the mass, and $x = 0$ is defined to be the equilibrium position. What is the force the spring exerts on the mass? Interpret the sign of the result.

▷ Differentiating, we find

$$\begin{aligned} F &= -\frac{dU}{dx} \\ &= -kx. \end{aligned}$$

If x is positive, then the force is negative, i.e., it acts so as to bring the mass back to equilibrium, and similarly for $x < 0$ we have $F > 0$.

Most books do the $F = -kx$ form before the $U = (1/2)kx^2$ form, and call it Hooke’s law. Neither form is really more fundamental than the other — we can always get from one to the other by integrating or differentiating.

Newton’s law of gravity

example 41

▷ Given the equation $U = -Gm_1 m_2 / r$ for the energy of gravitational interactions, find the corresponding equation for the gravitational force on mass m_2 . Interpret the positive and negative signs.

▷ We have to be a little careful here, because we’ve been taking r to be positive by definition, whereas the position, x , of mass m_2 could be positive or negative, depending on which side of m_1 it’s on.

For positive x , we have $r = x$, and differentiation gives

$$\begin{aligned} F &= -\frac{dU}{dx} \\ &= -Gm_1 m_2 / x^2. \end{aligned}$$

As in the preceding example, we have $F < 0$ when x is positive, because the object is being attracted back toward $x = 0$.

When x is negative, the relationship between r and x becomes $r = -x$, and the result for the force is the same as before, but with a minus sign. We can combine the two equations by writing

$$|F| = Gm_1m_2/r^2,$$

and this is the form traditionally known as Newton's law of gravity. As in the preceding example, the U and F equations contain equivalent information, and neither is more fundamental than the other.

Equilibrium

example 42

I previously described the condition for equilibrium as a local maximum or minimum of U . A differentiable function has a zero derivative at its extrema, and we can now relate this directly to force: zero force acts on an object when it is at equilibrium.

3.3 Resonance

Resonance is a phenomenon in which an oscillator responds most strongly to a driving force that matches its own natural frequency of vibration. For example, suppose a child is on a playground swing with a natural frequency of 1 Hz. That is, if you pull the child away from equilibrium, release her, and then stop doing anything for a while, she'll oscillate at 1 Hz. If there was no friction, as we assumed in section 2.5, then the sum of her gravitational and kinetic energy would remain constant, and the amplitude would be exactly the same from one oscillation to the next. However, friction is going to convert these forms of energy into heat, so her oscillations would gradually die out. To keep this from happening, you might give her a push once per cycle, i.e., the frequency of your pushes would be 1 Hz, which is the same as the swing's natural frequency. As long as you stay in rhythm, the swing responds quite well. If you start the swing from rest, and then give pushes at 1 Hz, the swing's amplitude rapidly builds up, as in figure a, until after a while it reaches a steady state in which friction removes just as much energy as you put in over the course of one cycle.

self-check F

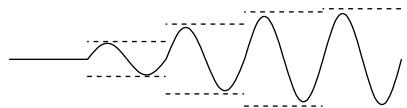
In figure a, compare the amplitude of the cycle immediately following the first push to the amplitude after the second. Compare the energies as well.

▷ Answer, p. 1059

What will happen if you try pushing at 2 Hz? Your first push puts in some momentum, p , but your second push happens after only half a cycle, when the swing is coming right back at you, with momentum $-p$! The momentum transfer from the second push is exactly enough to stop the swing. The result is a very weak, and not very sinusoidal, motion, b.

Making the math easy

This is a simple and physically transparent example of resonance: the swing responds most strongly if you match its natural rhythm. However, it has some characteristics that are mathematically ugly and possibly unrealistic. The quick, hard pushes are known as *impulse* forces, c, and they lead to an x - t graph that has nondifferentiable kinks. Impulsive forces like this are not only badly behaved mathematically, they are usually undesirable in practical terms. In a car engine, for example, the engineers work very hard to make the force on the pistons change smoothly, to avoid excessive vibration. Throughout the rest of this section, we'll assume a driving force that is sinusoidal, d, i.e., one whose F - t graph is either a sine function or a function that differs from a sine wave in phase, such as a cosine. The force is positive for half of each cycle and negative for the other half, i.e., there is both pushing and pulling. Sinusoidal functions have many nice mathematical characteristics (we can differentiate and integrate them, and the sum of sinusoidal functions



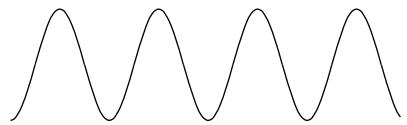
a / An x -versus- t graph for a swing pushed at resonance.



b / A swing pushed at twice its resonant frequency.



c / The F -versus- t graph for an impulsive driving force.



d / A sinusoidal driving force.

that have the same frequency is a sinusoidal function), and they are also used in many practical situations. For instance, my garage door zapper sends out a sinusoidal radio wave, and the receiver is tuned to resonance with it.

A second mathematical issue that I glossed over in the swing example was how friction behaves. In section 3.2.4, about forces between solids, the empirical equation for kinetic friction was independent of velocity. Fluid friction, on the other hand, is velocity-dependent. For a child on a swing, fluid friction is the most important form of friction, and is approximately proportional to v^2 . In still other situations, e.g., with a low-density gas or friction between solid surfaces that have been lubricated with a fluid such as oil, we may find that the frictional force has some other dependence on velocity, perhaps being proportional to v , or having some other complicated velocity dependence that can't even be expressed with a simple equation. It would be extremely complicated to have to treat all of these different possibilities in complete generality, so for the rest of this section, we'll assume friction proportional to velocity

$$F = -bv,$$

simply because the resulting equations happen to be the easiest to solve. Even when the friction doesn't behave in exactly this way, many of our results may still be at least qualitatively correct.

3.3.1 Damped, free motion

Numerical treatment

An oscillator that has friction is referred to as damped. Let's use numerical techniques to find the motion of a damped oscillator that is released away from equilibrium, but experiences no driving force after that. We can expect that the motion will consist of oscillations that gradually die out.

In section 2.5, we simulated the undamped case using our tried and true Python function based on conservation of energy. Now, however, that approach becomes a little awkward, because it involves splitting up the path to be traveled into n tiny segments, but in the presence of damping, each swing is a little shorter than the last one, and we don't know in advance exactly how far the oscillation will get before turning around. An easier technique here is to use force rather than energy. Newton's second law, $a = F/m$, gives $a = (-kx - bv)/m$, where we've made use of the result of example 40 for the force exerted by the spring. This becomes a little prettier if we rewrite it in the form

$$ma + bv + kx = 0,$$

which gives symmetric treatment to three terms involving x and its first and second derivatives, v and a . Now instead of calculating

the time $\Delta t = \Delta x/v$ required to move a predetermined distance Δx , we pick Δt and determine the distance traveled in that time, $\Delta x = v\Delta t$. Also, we can no longer update v based on conservation of energy, since we don't have any easy way to keep track of how much mechanical energy has been changed into heat energy. Instead, we recalculate the velocity using $\Delta v = a\Delta t$.

```

1 import math
2 k=39.4784          # chosen to give a period of 1 second
3 m=1.
4 b=0.211           # chosen to make the results simple
5 x=1.
6 v=0.
7 t=0.
8 dt=.01
9 n=1000
10 for j in range(n):
11     x=x+v*dt
12     a=(-k*x-b*v)/m
13     if (v>0) and (v+a*dt<0) :
14         print("turnaround at t=",t,", x=",x)
15     v=v+a*dt
16     t=t+dt

turnaround at t= 0.99 , x= 0.899919262445
turnaround at t= 1.99 , x= 0.809844934046
turnaround at t= 2.99 , x= 0.728777519477
turnaround at t= 3.99 , x= 0.655817260033
turnaround at t= 4.99 , x= 0.590154191135
turnaround at t= 5.99 , x= 0.531059189965
turnaround at t= 6.99 , x= 0.477875914756
turnaround at t= 7.99 , x= 0.430013546991
turnaround at t= 8.99 , x= 0.386940256644
turnaround at t= 9.99 , x= 0.348177318484

```

The spring constant, $k = 4\pi = 39.4784$ N/m, is designed so that if the undamped equation $f = (1/2\pi)\sqrt{k/m}$ was still true, the frequency would be 1 Hz. We start by noting that the addition of a small amount of damping doesn't seem to have changed the period at all, or at least not to within the accuracy of the calculation.¹⁰ You can check for yourself, however, that a large value of b , say 5 N·s/m, does change the period significantly.

We release the mass from $x = 1$ m, and after one cycle, it only comes back to about $x = 0.9$ m. I chose $b = 0.211$ N·s/m by fiddling

¹⁰This subroutine isn't as accurate a way of calculating the period as the energy-based one we used in the undamped case, since it only checks whether the mass turned around at some point during the time interval Δt .

around until I got this result, since a decrease of exactly 10% is easy to discuss. Notice how the amplitude after two cycles is about 0.81 m, i.e., 1 m times 0.9^2 : the amplitude has again dropped by exactly 10%. This pattern continues for as long as the simulation runs, e.g., for the last two cycles, we have $0.34818/0.38694=0.89982$, or almost exactly 0.9 again. It might have seemed capricious when I chose to use the unrealistic equation $F = -bv$, but this is the payoff. Only with $-bv$ friction do we get this kind of mathematically simple exponential decay.

Because the decay is exponential, it never dies out completely; this is different from the behavior we would have had with Coulomb friction, which does make objects grind completely to a stop at some point. With friction that acts like $F = -bv$, v gets smaller as the oscillations get smaller. The smaller and smaller force then causes them to die out at a rate that is slower and slower.

Analytic treatment

Taking advantage of this unexpectedly simple result, let's find an analytic solution for the motion. The numerical output suggests that we assume a solution of the form

$$x = Ae^{-ct} \sin(\omega_f t + \delta),$$

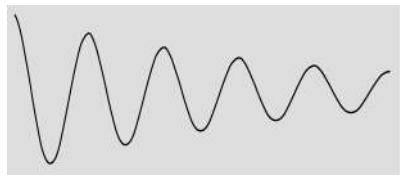
where the unknown constants ω_f and c will presumably be related to m , b , and k . The constant c indicates how quickly the oscillations die out. The constant ω_f is, as before, defined as 2π times the frequency, with the subscript f to indicate a free (undriven) solution. All our equations will come out much simpler if we use ω_s everywhere instead of ω_{fs} from now on, and, as physicists often do, I'll generally use the word "frequency" to refer to ω when the context makes it clear what I'm talking about. The phase angle δ has no real physical significance, since we can define $t = 0$ to be any moment in time we like.

self-check G

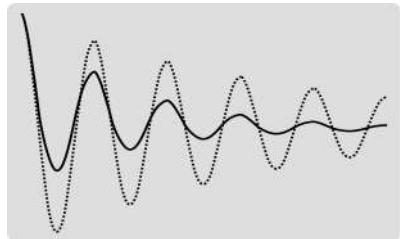
In figure f, which graph has the greater value of c ? \rightarrow Answer, p. 1059

The factor A for the initial amplitude can also be omitted without loss of generality, since the equation we're trying to solve, $ma + bv + kx = 0$, is linear. That is, v and a are the first and second derivatives of x , and the derivative of Ax is simply A times the derivative of x . Thus, if $x(t)$ is a solution of the equation, then multiplying it by a constant gives an equally valid solution. This is another place where we see that a damping force proportional to v is the easiest to handle mathematically. For a damping force proportional to v^2 , for example, we would have had to solve the equation $ma + bv^2 + kx = 0$, which is nonlinear.

For the purpose of determining ω_f and c , the most general form we need to consider is therefore $x = e^{-ct} \sin \omega_f t$, whose first and



e / A damped sine wave, of the form $x = Ae^{-ct} \sin(\omega_f t + \delta)$.



f / Self-check G.

second derivatives are $v = e^{-ct}(-c \sin \omega_f t + \omega \cos \omega_f t)$ and $a = e^{-ct}(c^2 \sin \omega_f t - 2\omega_f c \cos \omega_f t - \omega_f^2 \sin \omega_f t)$. Plugging these into the equation $ma + bv + kx = 0$ and setting the sine and cosine parts equal to zero gives, after some tedious algebra,

$$c = \frac{b}{2m}$$

and

$$\omega_f = \sqrt{\frac{k}{m} - \frac{b^2}{4m^2}}.$$

Intuitively, we expect friction to “slow down” the motion, as when we ride a bike into a big patch of mud. “Slow down,” however, could have more than one meaning here. It could mean that the oscillator would take more time to complete each cycle, or it could mean that as time went on, the oscillations would die out, thus giving smaller velocities.

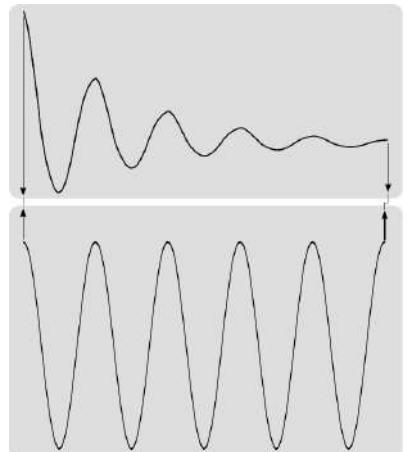
Our mathematical results show that both of these things happen. The first equation says that c , which indicates how quickly the oscillations damp out, is directly related to b , the strength of the damping.

The second equation, for the frequency, can be compared with the result from page 118 of $\sqrt{k/m}$ for the undamped system. Let’s refer to this now as ω_0 , to distinguish it from the actual frequency ω_f of the free oscillations when damping is present. The result for ω_f will be less than ω_0 , due to the presence of the $b^2/4m^2$ term. This tells us that the addition of friction to the system does increase the time required for each cycle. However, it is very common for the $b^2/4m^2$ term to be negligible, so that $\omega_f \approx \omega_0$.

Figure g shows an example. The damping here is quite strong: after only one cycle of oscillation, the amplitude has already been reduced by a factor of 2, corresponding to a factor of 4 in energy. However, the frequency of the damped oscillator is only about 1% lower than that of the undamped one; after five periods, the accumulated lag is just barely visible in the offsetting of the arrows. We can see that extremely strong damping — even stronger than this — would have been necessary in order to make $\omega_f \approx \omega_0$ a poor approximation.

3.3.2 The quality factor

It’s usually impractical to measure b directly and determine c from the equation $c = b/2m$. For a child on a swing, measuring b would require putting the child in a wind tunnel! It’s usually much easier to characterize the amount of damping by observing the actual damped oscillations and seeing how many cycles it takes for the mechanical energy to decrease by a certain factor. The unitless



g / A damped sine wave is compared with an undamped one, with m and k kept the same and only b changed.

quality factor, Q , is defined as $Q = \omega_0/2c$, and in the limit of weak damping, where $\omega \approx \omega_0$, this can be interpreted as the number of cycles required for the mechanical energy to fall off by a factor of $e^{2\pi} = 535.49\dots$. Using this new quantity, we can rewrite the equation for the frequency of damped oscillations in the slightly more elegant form $\omega_f = \omega_0\sqrt{1 - 1/4Q^2}$.

self-check H

What if we wanted to make a simpler definition of Q , as the number of oscillations required for the vibrations to die out completely, rather than the number required for the energy to fall off by this obscure factor? ▷ Answer, p. 1060

A graph

example 43

The damped motion in figure g has $Q \approx 4.5$, giving $\sqrt{1 - 1/4Q^2} \approx 0.99$, as claimed at the end of the preceding subsection.

Exponential decay in a trumpet

example 44

▷ The vibrations of the air column inside a trumpet have a Q of about 10. This means that even after the trumpet player stops blowing, the note will keep sounding for a short time. If the player suddenly stops blowing, how will the sound intensity 20 cycles later compare with the sound intensity while she was still blowing?

▷ The trumpet's Q is 10, so after 10 cycles the energy will have fallen off by a factor of 535. After another 10 cycles we lose another factor of 535, so the sound intensity is reduced by a factor of $535 \times 535 = 2.9 \times 10^5$.

The decay of a musical sound is part of what gives it its character, and a good musical instrument should have the right Q , but the Q that is considered desirable is different for different instruments. A guitar is meant to keep on sounding for a long time after a string has been plucked, and might have a Q of 1000 or 10000. One of the reasons why a cheap synthesizer sounds so bad is that the sound suddenly cuts off after a key is released.

3.3.3 Driven motion

The driven case is extremely important in science, technology, and engineering. We have an external driving force $F = F_m \sin \omega t$, where the constant F_m indicates the maximum strength of the force in either direction. The equation of motion is now

$$[1] \quad ma + bv + kx = F_m \sin \omega t$$

[equation of motion for a driven oscillator].

After the driving force has been applied for a while, we expect that the amplitude of the oscillations will approach some constant value. This motion is known as the *steady state*, and it's the most interesting thing to find out; as we'll see later, the most general type of motion is only a minor variation on the steady-state motion. For

Summary of Notation

k	spring constant
m	mass of the oscillator
b	sets the amount of damping, $F = -bv$
T	period
f	frequency, $1/T$
ω	(Greek letter omega), angular frequency, $2\pi f$, often referred to simply as "frequency"
ω_0	frequency the oscillator would have without damping, $\sqrt{k/m}$
ω_f	frequency of the free vibrations
c	sets the time scale for the exponential decay envelope e^{-ct} of the free vibrations
F_m	strength of the driving force, which is assumed to vary sinusoidally with frequency ω
A	amplitude of the steady-state response
δ	phase angle of the steady-state response

the steady-state motion, we're going to look for a solution of the form

$$x = A \sin(\omega t + \delta).$$

In contrast to the undriven case, here it's not possible to sweep A and δ under the rug. The amplitude of the steady-state motion, A , is actually the most interesting thing to know about the steady-state motion, and it's not true that we still have a solution no matter how we fiddle with A ; if we have a solution for a certain value of A , then multiplying A by some constant would break the equality between the two sides of the equation of motion. It's also no longer true that we can get rid of δ simply by redefining when we start the clock; here δ represents a *difference* in time between the start of one cycle of the driving force and the start of the corresponding cycle of the motion.

The velocity and acceleration are $v = \omega A \cos(\omega t + \delta)$ and $a = -\omega^2 A \sin(\omega t + \delta)$, and if we plug these into the equation of motion, [1], and simplify a little, we find

$$[2] \quad (k - m\omega^2) \sin(\omega t + \delta) + \omega b \cos(\omega t + \delta) = \frac{F_m}{A} \sin \omega t.$$

The sum of any two sinusoidal functions with the same frequency is also a sinusoidal, so the whole left side adds up to a sinusoidal. By fiddling with A and δ we can make the amplitudes and phases of the two sides of the equation match up.

Steady state, no damping

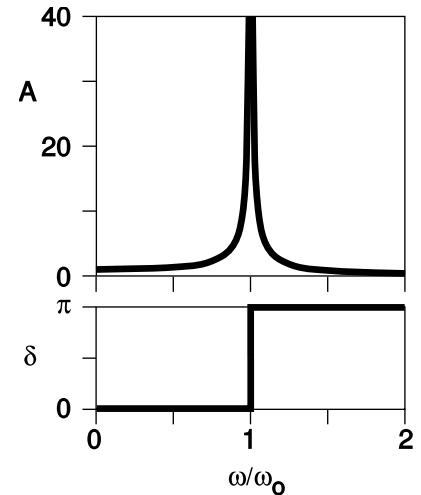
A and δ are easy to find in the case where there is no damping at all. There are now no cosines in equation [2] above, only sines, so if we wish we can set δ to zero, and we find $A = F_m/(k - m\omega^2) = F_m/m(\omega_0^2 - \omega^2)$. This, however, makes A negative for $\omega > \omega_0$. The variable δ was designed to represent this kind of phase relationship, so we prefer to keep A positive and set $\delta = \pi$ for $\omega > \omega_0$. Our results are then

$$A = \frac{F_m}{m|\omega^2 - \omega_0^2|}$$

and

$$\delta = \begin{cases} 0, & \omega < \omega_0 \\ \pi, & \omega > \omega_0 \end{cases}.$$

The most important feature of the result is that there is a resonance: the amplitude becomes greater and greater, and approaches infinity, as ω approaches the resonant frequency ω_0 . This is the physical behavior we anticipated on page 175 in the example of pushing a child on a swing. If the driving frequency matches the frequency of the free vibrations, then the driving force will always be in the



h / Dependence of the amplitude and phase angle on the driving frequency, for an undamped oscillator. The amplitudes were calculated with F_m , m , and ω_0 , all set to 1.

right direction to add energy to the swing. At a driving frequency very different from the resonant frequency, we might get lucky and push at the right time during one cycle, but our next push would come at some random point in the next cycle, possibly having the effect of slowing the swing down rather than speeding it up.

The interpretation of the infinite amplitude at $\omega = \omega_0$ is that there really isn't any steady state if we drive the system exactly at resonance — the amplitude will just keep on increasing indefinitely. In real life, the amplitude can't be infinite both because there is always some damping and because there will always be some difference, however small, between ω and ω_0 . Even though the infinity is unphysical, it has entered into the popular consciousness, starting with the eccentric Serbian-American inventor and physicist Nikola Tesla. Around 1912, the tabloid newspaper *The World Today* credulously reported a story which Tesla probably fabricated — or wildly exaggerated — for the sake of publicity. Supposedly he created a steam-powered device “no larger than an alarm clock,” containing a piston that could be made to vibrate at a tunable and precisely controlled frequency. “He put his little vibrator in his coat-pocket and went out to hunt a half-erected steel building. Down in the Wall Street district, he found one — ten stories of steel framework without a brick or a stone laid around it. He clamped the vibrator to one of the beams, and fussed with the adjustment [presumably hunting for the building’s resonant frequency] until he got it. Tesla said finally the structure began to creak and weave and the steel-workers came to the ground panic-stricken, believing that there had been an earthquake. Police were called out. Tesla put the vibrator in his pocket and went away. Ten minutes more and he could have laid the building in the street. And, with the same vibrator he could have dropped the Brooklyn Bridge into the East River in less than an hour.”

The phase angle δ also exhibits surprising behavior. As the frequency is tuned upward past resonance, the phase abruptly shifts so that the phase of the response is opposite to that of the driving force. There is a simple interpretation for this. The system’s mechanical energy can only change due to work done by the driving force, since there is no damping to convert mechanical energy to heat. In the steady state, then, the power transmitted by the driving force over a full cycle of motion must average out to zero. In general, the work theorem $dE = F dx$ can always be divided by dt on both sides to give the useful relation $P = Fv$. If Fv is to average out to zero, then F and v must be out of phase by $\pm\pi/2$, and since v is ahead of x by a phase angle of $\pi/2$, the phase angle between x and F must be zero or π .

Given that these are the two possible phases, why is there a difference in behavior between $\omega < \omega_0$ and $\omega > \omega_0$? At the low-frequency limit, consider $\omega = 0$, i.e., a constant force. A constant

force will simply displace the oscillator to one side, reaching an equilibrium that is offset from the usual one. The force and the response are in phase, e.g., if the force is to the right, the equilibrium will be offset to the right. This is the situation depicted in the amplitude graph of figure h at $\omega = 0$. The response, which is not zero, is simply this static displacement of the oscillator to one side.

At high frequencies, on the other hand, imagine shaking the poor child on the swing back and forth with a force that oscillates at 10 Hz. This is so fast that there is essentially no time for the force $F = -kx$ from gravity and the chain to act from one cycle to the next. The problem becomes equivalent to the oscillation of a *free* object. If the driving force varies like $\sin(\omega t)$, with $\delta = 0$, then the acceleration is also proportional to the sine. Integrating, we find that the velocity goes like minus a cosine, and a second integration gives a position that varies as minus the sine — opposite in phase to the driving force. Intuitively, this mathematical result corresponds to the fact that at the moment when the object has reached its maximum displacement to the *right*, that is the time when the greatest force is being applied to the *left*, in order to turn it around and bring it back toward the center.

A practice mute for a violin

example 45

The amplitude of the driven vibrations, $A = F_m/(m|\omega^2 - \omega_0^2|)$, contains an inverse proportionality to the mass of the vibrating object. This is simply because a given force will produce less acceleration when applied to a more massive object. An application is shown in figure 45.

In a stringed instrument, the strings themselves don't have enough surface area to excite sound waves very efficiently. In instruments of the violin family, as the strings vibrate from left to right, they cause the bridge (the piece of wood they pass over) to wiggle clockwise and counterclockwise, and this motion is transmitted to the top panel of the instrument, which vibrates and creates sound waves in the air.

A string player who wants to practice at night without bothering the neighbors can add some mass to the bridge. Adding mass to the bridge causes the amplitude of the vibrations to be smaller, and the sound to be much softer. A similar effect is seen when an electric guitar is used without an amp. The body of an electric guitar is so much more massive than the body of an acoustic guitar that the amplitude of its vibrations is very small.

i / Example 45: a viola without a mute (left), and with a mute (right). The mute doesn't touch the strings themselves.



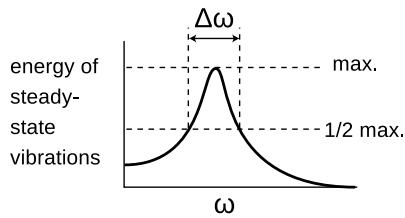
Steady state, with damping

The extension of the analysis to the damped case involves some lengthy algebra, which I've outlined on page 1027 in appendix 2. The results are shown in figure j. It's not surprising that the steady state response is weaker when there is more damping, since the steady state is reached when the power extracted by damping matches the power input by the driving force. The maximum amplitude, at the peak of the resonance curve, is approximately proportional to Q .

self-check i

From the final result of the analysis on page 1027, substitute $\omega = \omega_0$, and satisfy yourself that the result is proportional to Q . Why is $A_{res} \propto Q$ only an approximation?

▷ Answer, p. 1060



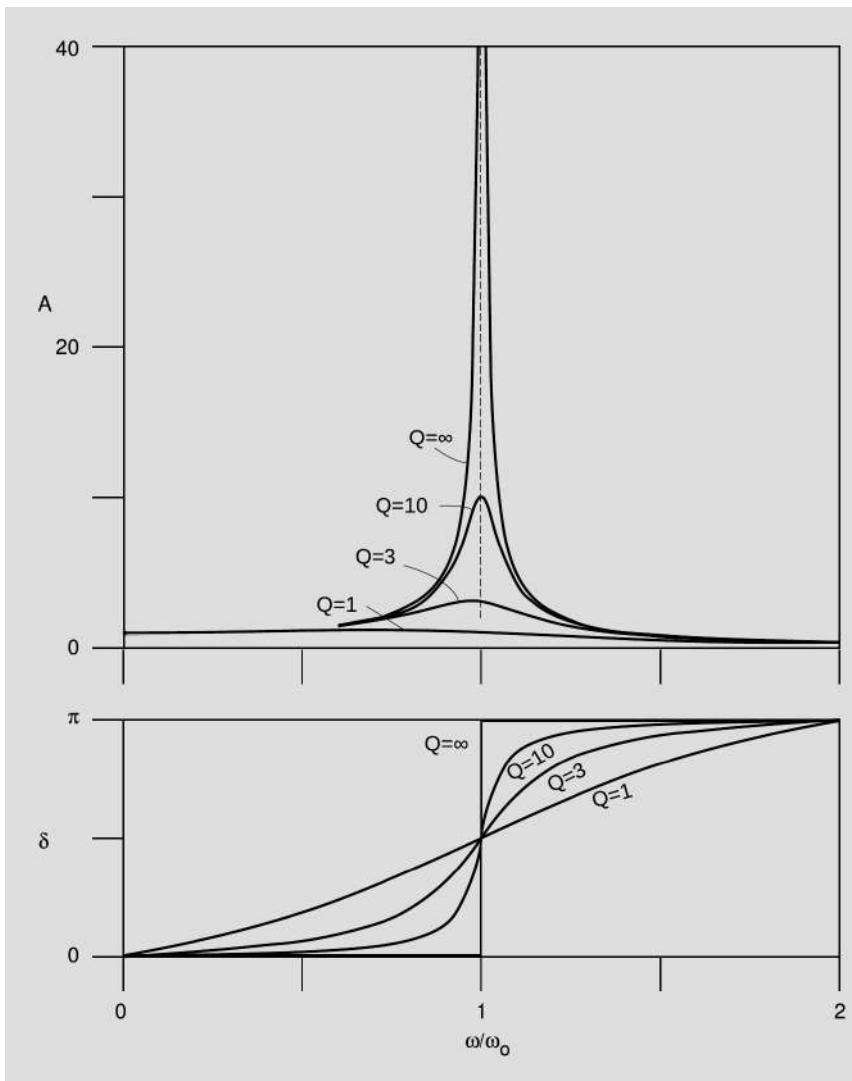
k / The definition of $\Delta\omega$, the full width at half maximum.

What is surprising is that the amplitude is strongly affected by damping close to resonance, but only weakly affected far from it. In other words, the shape of the resonance curve is broader with more damping, and even if we were to scale up a high-damping curve so that its maximum was the same as that of a low-damping curve, it would still have a different shape. The standard way of describing the shape numerically is to give the quantity $\Delta\omega$, called the *full width at half-maximum*, or FWHM, which is defined in figure k. Note that the y axis is energy, which is proportional to the square of the amplitude. Our previous observations amount to a statement that $\Delta\omega$ is greater when the damping is stronger, i.e., when the Q is lower. It's not hard to show from the equations on page 1027 that for large Q , the FWHM is given approximately by

$$\Delta\omega \approx \omega_0/Q.$$

Another thing we notice in figure j is that for small values of Q the frequency ω_{res} of the maximum A is less than ω_0 .¹¹ At even

¹¹The relationship is $\omega_{max\ A}/\omega_0 = \sqrt{1 - 1/2Q^2}$, which is similar in form to the equation for the frequency of the free vibration, $\omega_f/\omega_0 = \sqrt{1 - 1/4Q^2}$. A subtle point here is that although the maximum of A and the maximum of A^2 must occur at the same frequency, the maximum energy does not occur, as we might expect, at the same frequency as the maximum of A^2 . This is because



j / Dependence of the amplitude and phase angle on the driving frequency. The undamped case is $Q = \infty$, and the other curves represent $Q=1$, 3, and 10. F_m , m , and ω_0 are all set to 1.

lower values of Q , like $Q = 1$, the $A - \omega$ curve doesn't even have a maximum near $\omega > 0$.

An opera singer breaking a wineglass

example 46

In order to break a wineglass by singing, an opera singer must first tap the glass to find its natural frequency of vibration, and then sing the same note back, so that her driving force will produce a response with the greatest possible amplitude. If she's shopping for the right glass to use for this display of her prowess, she should look for one that has the greatest possible Q , since the resonance curve has a higher maximum for higher values of Q .

the interaction energy is proportional to A^2 regardless of frequency, but the kinetic energy is proportional to $A^2\omega^2$. The maximum energy actually occurs are precisely ω_0 .

Collapse of the Nimitz Freeway

example 47

Figure I shows a section of the Nimitz Freeway in Oakland, CA, that collapsed during an earthquake in 1989. An earthquake consists of many low-frequency vibrations that occur simultaneously, which is why it sounds like a rumble of indeterminate pitch rather than a low hum. The frequencies that we can hear are not even the strongest ones; most of the energy is in the form of vibrations in the range of frequencies from about 1 Hz to 10 Hz.



I / The collapsed section of the Nimitz Freeway

All the structures we build are resting on geological layers of dirt, mud, sand, or rock. When an earthquake wave comes along, the topmost layer acts like a system with a certain natural frequency of vibration, sort of like a cube of jello on a plate being shaken from side to side. The resonant frequency of the layer depends on how stiff it is and also on how deep it is. The ill-fated section of the Nimitz freeway was built on a layer of mud, and analysis by geologist Susan E. Hough of the U.S. Geological Survey shows that the mud layer's resonance was centered on about 2.5 Hz, and had a width covering a range from about 1 Hz to 4 Hz.

When the earthquake wave came along with its mixture of frequencies, the mud responded strongly to those that were close to its own natural 2.5 Hz frequency. Unfortunately, an engineering analysis after the quake showed that the overpass itself had a resonant frequency of 2.5 Hz as well! The mud responded strongly to the earthquake waves with frequencies close to 2.5 Hz, and the bridge responded strongly to the 2.5 Hz vibrations of the mud, causing sections of it to collapse.

Physical reason for the relationship between Q and the FWHM

What is the reason for this surprising relationship between the damping and the width of the resonance? Fundamentally, it has to do with the fact that friction causes a system to lose its “memory” of its previous state. If the Pioneer 10 space probe, coasting through the frictionless vacuum of interplanetary space, is detected by aliens a million years from now, they will be able to trace its trajectory backwards and infer that it came from our solar system. On the other hand, imagine that I shove a book along a tabletop, it comes to rest, and then someone else walks into the room. There will be no clue as to which direction the book was moving before it stopped — friction has erased its memory of its motion. Now consider the playground swing driven at twice its natural frequency, figure m, where the undamped case is repeated from figure b on page 175. In the undamped case, the first push starts the swing moving with momentum p , but when the second push comes, if there is no friction at all, it now has a momentum of exactly $-p$, and the momentum transfer from the second push is exactly enough to stop it dead. With moderate damping, however, the momentum on the rebound is not quite $-p$, and the second push’s effect isn’t quite as disas-

trous. With very strong damping, the swing comes essentially to rest long before the second push. It has lost all its memory, and the second push puts energy into the system rather than taking it out. Although the detailed mathematical results with this kind of impulsive driving force are different,¹² the general results are the same as for sinusoidal driving: the less damping there is, the greater the penalty you pay for driving the system off of resonance.

High-Q speakers

example 48

Most good audio speakers have $Q \approx 1$, but the resonance curve for a higher- Q oscillator always lies above the corresponding curve for one with a lower Q , so people who want their car stereos to be able to rattle the windows of the neighboring cars will often choose speakers that have a high Q . Of course they could just use speakers with stronger driving magnets to increase F_m , but the speakers might be more expensive, and a high- Q speaker also has less friction, so it wastes less energy as heat.

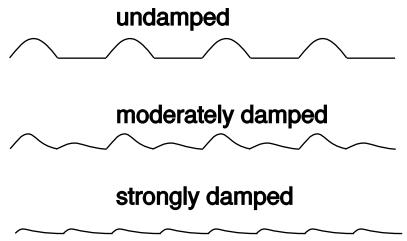
One problem with this is that whereas the resonance curve of a low- Q speaker (its “response curve” or “frequency response” in audiophile lingo) is fairly flat, a higher- Q speaker tends to emphasize the frequencies that are close to its natural resonance. In audio, a flat response curve gives more realistic reproduction of sound, so a higher quality factor, Q , really corresponds to a *lower-quality* speaker.

Another problem with high- Q speakers is discussed in example 51 on page 189 .

Changing the pitch of a wind instrument

example 49

- ▷ A saxophone player normally selects which note to play by choosing a certain fingering, which gives the saxophone a certain resonant frequency. The musician can also, however, change the pitch significantly by altering the tightness of her lips. This corresponds to driving the horn slightly off of resonance. If the pitch can be altered by about 5% up or down (about one musical half-step) without too much effort, roughly what is the Q of a saxophone?
- ▷ Five percent is the width on one side of the resonance, so the full width is about 10%, $\Delta f/f_0 \approx 0.1$. The equation $\Delta\omega = \omega_0/Q$ is defined in terms of angular frequency, $\omega = 2\pi f$, and we’ve been given our data in terms of ordinary frequency, f . The factors of 2π



m / An x -versus- t graph of the steady-state motion of a swing being pushed at twice its resonant frequency by an impulsive force.

¹²For example, the graphs calculated for sinusoidal driving have resonances that are somewhat below the natural frequency, getting lower with increasing damping, until for $Q \leq 1$ the maximum response occurs at $\omega = 0$. In figure m, however, we can see that impulsive driving at $\omega = 2\omega_0$ produces a steady state with more energy than at $\omega = \omega_0$.

end up canceling out, however:

$$\begin{aligned} Q &= \frac{\omega_0}{\Delta\omega} \\ &= \frac{2\pi f_0}{2\pi\Delta f} \\ &= \frac{f_0}{f} \\ &\approx 10 \end{aligned}$$

In other words, once the musician stops blowing, the horn will continue sounding for about 10 cycles before its energy falls off by a factor of 535. (Blues and jazz saxophone players will typically choose a mouthpiece that gives a low Q , so that they can produce the bluesy pitch-slides typical of their style. “Legit,” i.e., classically oriented players, use a higher- Q setup because their style only calls for enough pitch variation to produce a vibrato, and the higher Q makes it easier to play in tune.)

Q of a radio receiver *example 50*

- ▷ A radio receiver used in the FM band needs to be tuned in to within about 0.1 MHz for signals at about 100 MHz. What is its Q ?
- ▷ As in the last example, we’re given data in terms of f_s , not ω_s , but the factors of 2π cancel. The resulting Q is about 1000, which is extremely high compared to the Q values of most mechanical systems.

Transients

What about the motion before the steady state is achieved? When we computed the undriven motion numerically on page 176, the program had to initialize the position and velocity. By changing these two variables, we could have gotten any of an infinite number of simulations.¹³ The same is true when we have an equation of motion with a driving term, $ma + bv + kx = F_m \sin \omega t$ (p. 180, equation [1]). The steady-state solutions, however, have no adjustable parameters at all — A and δ are uniquely determined by the parameters of the driving force and the oscillator itself. If the oscillator isn’t initially in the steady state, then it will not have the steady-state motion at first. What kind of motion will it have?

The answer comes from realizing that if we start with the solution to the driven equation of motion, and then add to it any solution to the free equation of motion, the result,

$$x = A \sin(\omega t + \delta) + A' e^{-ct} \sin(\omega_f t + \delta'),$$

¹³If you’ve learned about differential equations, you’ll know that any second-order differential equation requires the specification of two boundary conditions in order to specify solution uniquely.

is also a solution of the driven equation. Here, as before, ω_f is the frequency of the free oscillations ($\omega_f \approx \omega_0$ for small Q), ω is the frequency of the driving force, A and δ are related as usual to the parameters of the driving force, and A' and δ' can have any values at all. Given the initial position and velocity, we can always choose A' and δ' to reproduce them, but this is not something one often has to do in real life. What's more important is to realize that the second term dies out exponentially over time, decaying at the same rate at which a free vibration would. For this reason, the A' term is called a transient. A high- Q oscillator's transients take a long time to die out, while a low- Q oscillator always settles down to its steady state very quickly.

Boomy bass

example 51

In example 48 on page 187, I've already discussed one of the drawbacks of a high- Q speaker, which is an uneven response curve. Another problem is that in a high- Q speaker, transients take a long time to die out. The bleeding-eardrums crowd tend to focus mostly on making their bass loud, so it's usually their woofers that have high Q s. The result is that bass notes, "ring" after the onset of the note, a phenomenon referred to as "boomy bass."

Overdamped motion

The treatment of free, damped motion on page 178 skipped over a subtle point: in the equation $\omega_f = \sqrt{k/m - b^2/4m^2} = \omega_0\sqrt{1 - 1/4Q^2}$, $Q < 1/2$ results in an answer that is the square root of a negative number. For example, suppose we had $k = 0$, which corresponds to a neutral equilibrium. A physical example would be a mass sitting in a tub of syrup. If we set it in motion, it won't oscillate — it will simply slow to a stop. This system has $Q = 0$. The equation of motion in this case is $ma + bv = 0$, or, more suggestively,

$$m \frac{dv}{dt} + bv = 0.$$

One can easily verify that this has the solution $v = (\text{constant})e^{-bt/m}$, and integrating, we find $x = (\text{constant})e^{-bt/m} + (\text{constant})$. In other words, the reason ω_f comes out to be mathematical nonsense¹⁴ is that we were incorrect in assuming a solution that oscillated at a frequency ω_f . The actual motion is not oscillatory at all.

In general, systems with $Q < 1/2$, called overdamped systems, do not display oscillatory motion. Most cars' shock absorbers are designed with $Q \approx 1/2$, since it's undesirable for the car to undulate up and down for a while after you go over a bump. (Shocks with extremely low values of Q are not good either, because such a system takes a very long time to come back to equilibrium.) It's not par-

¹⁴Actually, if you know about complex numbers and Euler's theorem, it's not quite so nonsensical.

ticularly important for our purposes, but for completeness I'll note, as you can easily verify, that the general solution to the equation of motion for $0 < Q < 1/2$ is of the form $x = Ae^{-ct} + Be^{-dt}$, while $Q = 1/2$, called the critically damped case, gives $x = (A + Bt)e^{-ct}$.

3.4 Motion in three dimensions

3.4.1 The Cartesian perspective

When my friends and I were bored in high school, we used to play a paper-and-pencil game which, although we never knew it, was Very Educational — in fact, it pretty much embodies the entire worldview of classical physics. To play the game, you draw a racetrack on graph paper, and try to get your car around the track before anyone else. The default is for your car to continue at constant speed in a straight line, so if it moved three squares to the right and one square up on your last turn, it will do the same this turn. You can also control the car's motion by changing its Δx and Δy by up to one unit. If it moved three squares to the right last turn, you can have it move anywhere from two to four squares to the right this turn.

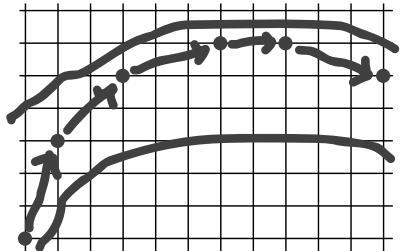
b / French mathematician René Descartes invented analytic geometry; Cartesian (xyz) coordinates are named after him. He did work in philosophy, and was particularly interested in the mind-body problem. He was a skeptic and an antistaristotelian, and, probably for fear of religious persecution, spent his adult life in the Netherlands, where he fathered a daughter with a Protestant peasant whom he could not marry. He kept his daughter's existence secret from his enemies in France to avoid giving them ammunition, but he was crushed when she died of scarlatina at age 5. A pious Catholic, he was widely expected to be sainted. His body was buried in Sweden but then reburied several times in France, and along the way everything but a few fingerbones was stolen by peasants who expected the body parts to become holy relics.

The fundamental way of dealing with the direction of an object's motion in physics is to use conservation of momentum, since momentum depends on direction. Up until now, we've only done momentum in one dimension. How does this relate to the racetrack game? In the game, the motion of a car from one turn to the next is represented by its Δx and Δy . In one dimension, we would only need Δx , which could be related to the velocity, $\Delta x/\Delta t$, and the momentum, $m\Delta x/\Delta t$. In two dimensions, the rules of the game amount to a statement that if there is no momentum transfer, then both $m\Delta x/\Delta t$ and $m\Delta y/\Delta t$ stay the same. In other words, there are two flavors of momentum, and they are *separately* conserved. All of this so far has been done with an artificial division of time into "turns," but we can fix that by redefining everything in terms of derivatives, and for motion in three dimensions rather than two, we augment x and y with z :

$$v_x = \frac{dx}{dt} \quad v_y = \frac{dy}{dt} \quad v_z = \frac{dz}{dt}$$

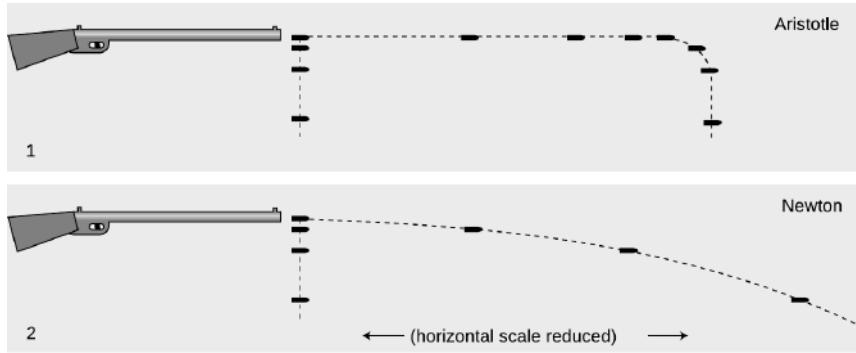
and

$$p_x = mv_x \quad p_y = mv_y \quad p_z = mv_z$$



a / The car can change its x and y motions by one square every turn.





c / Bullets are dropped and shot at the same time.

We call these the x , y , and z components of the velocity and the momentum.

There is both experimental and theoretical evidence that the x , y , and z momentum components are separately conserved, and that a momentum transfer (force) along one axis has no effect on the momentum components along the other two axes. On page 89, for example, I argued that it was impossible for an air hockey puck to make a 180-degree turn spontaneously, because then in the frame moving along with the puck, it would have begun moving after starting from rest. Now that we're working in two dimensions, we might wonder whether the puck could spontaneously make a 90-degree turn, but exactly the same line of reasoning shows that this would be impossible as well, which proves that the puck can't trade x -momentum for y -momentum. A more general proof of separate conservation will be given on page 218, after some of the appropriate mathematical techniques have been introduced.

As an example of the experimental evidence for separate conservation of the momentum components, figure c shows correct and incorrect predictions of what happens if you shoot a rifle and arrange for a second bullet to be dropped from the same height at exactly the same moment when the first one left the barrel. Nearly everyone expects that the dropped bullet will reach the dirt first, and Aristotle would have agreed, since he believed that the bullet had to lose its horizontal motion before it could start moving vertically. In reality, we find that the vertical momentum transfer between the earth and the bullet is completely unrelated to the horizontal momentum. The bullet ends up with $p_y < 0$, while the planet picks up an upward momentum $p_y > 0$, and the total momentum in the y direction remains zero. Both bullets hit the ground at the same time. This is much simpler than the Aristotelian version!

The Pelton waterwheel

example 52

▷ There is a general class of machines that either do work on a

gas or liquid, like a boat's propeller, or have work done on them by a gas or liquid, like the turbine in a hydroelectric power plant. Figure d shows two types of surfaces that could be attached to the circumference of an old-fashioned waterwheel. Compare the force exerted by the water in the two cases.

Let the x axis point to the right, and the y axis up. In both cases, the stream of water rushes down onto the surface with momentum $p_{y,i} = -p_0$, where the subscript i stands for "initial," i.e., before the collision.

In the case of surface 1, the streams of water leaving the surface have no momentum in the y direction, and their momenta in the x direction cancel. The final momentum of the water is zero along both axes, so its entire momentum, $-p_0$, has been transferred to the waterwheel.

When the water leaves surface 2, however, its momentum isn't zero. If we assume there is no friction, it's $p_{y,f} = +p_0$, with the positive sign indicating upward momentum. The change in the water's momentum is $p_{y,f} - p_{y,i} = 2p_0$, and the momentum transferred to the waterwheel is $-2p_0$.

Force is defined as the rate of transfer of momentum, so surface 2 experiences double the force. A waterwheel constructed in this way is known as a Pelton waterwheel.

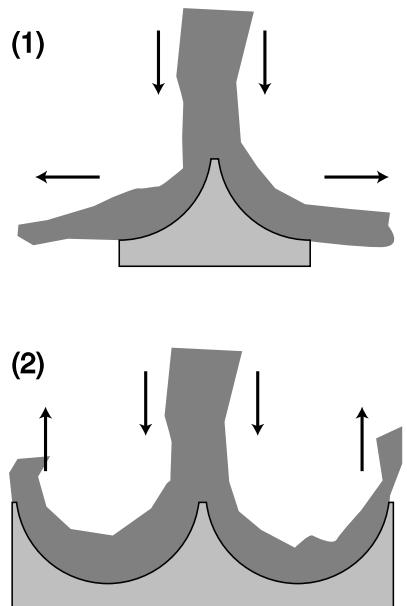
The Yarkovsky effect

example 53

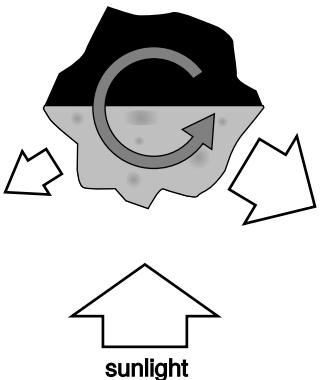
We think of the planets and asteroids as inhabiting their orbits permanently, but it is possible for an orbit to change over periods of millions or billions of years, due to a variety of effects. For asteroids with diameters of a few meters or less, an important mechanism is the Yarkovsky effect, which is easiest to understand if we consider an asteroid spinning about an axis that is exactly perpendicular to its orbital plane.

The illuminated side of the asteroid is relatively hot, and radiates more infrared light than the dark (night) side. Light has momentum, and a total force away from the sun is produced by combined effect of the sunlight hitting the asteroid and the imbalance between the momentum radiated away on the two sides. This force, however, doesn't cause the asteroid's orbit to change over time, since it simply cancels a tiny fraction of the sun's gravitational attraction. The result is merely a tiny, undetectable violation of Kepler's law of periods.

Consider the sideways momentum transfers, however. In figure e, the part of the asteroid on the right has been illuminated for half a spin-period (half a "day") by the sun, and is hot. It radiates more light than the morning side on the left. This imbalance produces a total force in the x direction which points to the left. If the asteroid's orbital motion is to the left, then this is a force in the same



d / Two surfaces that could be used to extract energy from a stream of water.



e / An asteroid absorbs visible light from the sun, and gets rid of the energy by radiating infrared light.

direction as the motion, which will do positive work, increasing the asteroid's energy and boosting it into an orbit with a greater radius. On the other hand, if the asteroid's spin and orbital motion are in opposite directions, the Yarkovsky push brings the asteroid spiraling in closer to the sun.

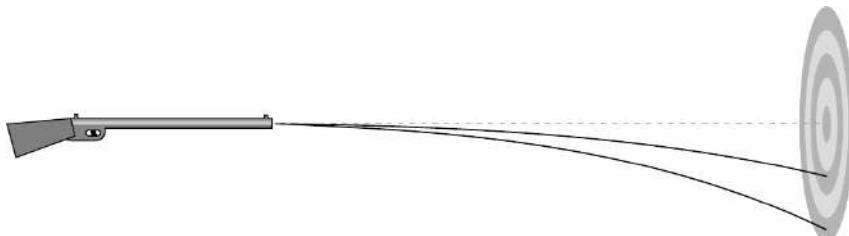
Calculations show that it takes on the order of 10^7 to 10^8 years for the Yarkovsky effect to move an asteroid out of the asteroid belt and into the vicinity of earth's orbit, and this is about the same as the typical age of a meteorite as estimated by its exposure to cosmic rays. The Yarkovsky effect doesn't remove all the asteroids from the asteroid belt, because many of them have orbits that are stabilized by gravitational interactions with Jupiter. However, when collisions occur, the fragments can end up in orbits which are not stabilized in this way, and they may then end up reaching the earth due to the Yarkovsky effect. The cosmic-ray technique is really telling us how long it has been since the fragment was broken out of its parent.

Discussion Questions

A The following is an incorrect explanation of a fact about target shooting:

"Shooting a high-powered rifle with a high muzzle velocity is different from shooting a less powerful gun. With a less powerful gun, you have to aim quite a bit above your target, but with a more powerful one you don't have to aim so high because the bullet doesn't drop as fast."

What is the correct explanation?



f / Discussion question A.

B You have thrown a rock, and it is flying through the air in an arc. If the earth's gravitational force on it is always straight down, why doesn't it just go straight down once it leaves your hand?

C Consider the example of the bullet that is dropped at the same moment another bullet is fired from a gun. What would the motion of the two bullets look like to a jet pilot flying alongside in the same direction as the shot bullet and at the same horizontal speed?

3.4.2 Rotational invariance

The Cartesian approach requires that we choose x , y , and z axes. How do we choose them correctly? The answer is that it had better not matter which directions the axes point (provided they're perpendicular to each other), or where we put the origin, because if it did matter, it would mean that space was asymmetric. If there was a certain point in the universe that was the right place to put the origin, where would it be? The top of Mount Olympus? The United Nations headquarters? We find that experiments come out the same no matter where we do them, and regardless of which way the laboratory is oriented, which indicates that no location in space or direction in space is special in any way.¹⁵

This is closely related to the idea of Galilean relativity stated on page 62, from which we already know that the absolute *motion* of a frame of reference is irrelevant and undetectable. Observers using frames of reference that are in motion relative to each other will not even agree on the permanent identity of a particular point in space, so it's not possible for the laws of physics to depend on where you are in space. For instance, if gravitational energies were proportional to $m_1 m_2$ in one location but to $(m_1 m_2)^{1.00001}$ in another, then it would be possible to determine when you were in a state of absolute motion, because the behavior of gravitational interactions would change as you moved from one region to the other.

Because of this close relationship, we restate the principle of Galilean relativity in a more general form. This extended principle of Galilean relativity states that the laws of physics are no different in one time and place than in another, and that they also don't depend on your orientation or your motion, provided that your motion is in a straight line and at constant speed.

The irrelevance of time and place could have been stated in chapter 1, but since this section is the first one in which we're dealing with three-dimensional physics in full generality, the irrelevance of orientation is what we really care about right now. This property of the laws of physics is called rotational invariance. The word "invariance" means a lack of change, i.e., the laws of physics don't change when we reorient our frame of reference.

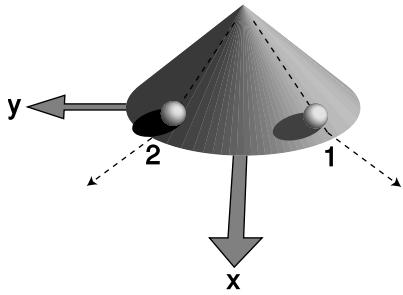
Rotational invariance of gravitational interactions example 54
Gravitational energies depend on the quantity $1/r$, which by the

¹⁵Of course, you could tell in a sealed laboratory which way was down, but that's because there happens to be a big planet nearby, and the planet's gravitational field reaches into the lab, not because space itself has a special down direction. Similarly, if your experiment was sensitive to magnetic fields, it might matter which way the building was oriented, but that's because the earth makes a magnetic field, not because space itself comes equipped with a north direction.

Pythagorean theorem equals

$$\frac{1}{\sqrt{\Delta x^2 + \Delta y^2 + \Delta z^2}}.$$

Rotating a line segment doesn't change its length, so this expression comes out the same regardless of which way we orient our coordinate axes. Even though Δx , Δy , and Δz are different in differently oriented coordinate systems, r is the same.



g / Two balls roll down a cone and onto a plane.

Kinetic energy

example 55

Kinetic energy equals $(1/2)mv^2$, but what does that mean in three dimensions, where we have v_x , v_y , and v_z ? If you were tempted to add the components and calculate $K = (1/2)m(v_x + v_y + v_z)^2$, figure g should convince you otherwise. Using that method, we'd have to assign a kinetic energy of zero to ball number 1, since its negative v_y would exactly cancel its positive v_x , whereas ball number 2's kinetic energy wouldn't be zero. This would violate rotational invariance, since the balls would behave differently.

The only possible way to generalize kinetic energy to three dimensions, without violating rotational invariance, is to use an expression that resembles the Pythagorean theorem,

$$v = \sqrt{v_x^2 + v_y^2 + v_z^2},$$

which results in

$$K = \frac{1}{2}m(v_x^2 + v_y^2 + v_z^2).$$

Since the velocity components are squared, the positive and negative signs don't matter, and the two balls in the example behave the same way.

3.4.3 Vectors

Remember the title of this book? It would have been possible to obtain the result of example 55 by applying the Pythagorean theorem to dx , dy , and dz , and then dividing by dt , but the rotational invariance approach is *simpler*, and is useful in a much broader context. Even with a quantity you presently know nothing about, say the magnetic field, you can infer that if the components of the magnetic field are B_x , B_y , and B_z , then the physically useful way to talk about the strength of the magnetic field is to define it as $\sqrt{B_x^2 + B_y^2 + B_z^2}$. Nature knows your brain cells are precious, and doesn't want you to have to waste them by memorizing mathematical rules that are different for magnetic fields than for velocities.

When mathematicians see that the same set of techniques is useful in many different contexts, that's when they start making definitions that allow them to stop reinventing the wheel. The ancient Greeks, for example, had no general concept of fractions. They couldn't say that a circle's radius divided by its diameter was equal to the number $1/2$. They had to say that the radius and the diameter were in the ratio of one to two. With this limited number concept, they couldn't have said that water was dripping out of a tank at a rate of $3/4$ of a barrel per day; instead, they would have had to say that over four days, three barrels worth of water would be lost. Once enough of these situations came up, some clever mathematician finally realized that it would make sense to define something called a fraction, and that one could think of these fraction thingies as numbers that lay in the gaps between the traditionally recognized numbers like zero, one, and two. Later generations of mathematicians introduced further subversive generalizations of the number concepts, inventing mathematical creatures like negative numbers, and the square root of two, which can't be expressed as a fraction.

In this spirit, we define a *vector* as any quantity that has both an amount and a direction in space. In contradistinction, a *scalar* has an amount, but no direction. Time and temperature are scalars. Velocity, acceleration, momentum, and force are vectors. In one dimension, there are only two possible directions, and we can use positive and negative numbers to indicate the two directions. In more than one dimension, there are infinitely many possible directions, so we can't use the two symbols $+$ and $-$ to indicate the direction of a vector. Instead, we can specify the three components of the vector, each of which can be either negative or positive. We represent vector quantities in handwriting by writing an arrow above them, so for example the momentum vector looks like this, \vec{p} , but the arrow looks ugly in print, so in books vectors are usually shown in bold-face type: \mathbf{p} . A straightforward way of thinking about vectors is that a vector equation really represents three different equations. For instance, conservation of momentum could be written in terms

of the three components,

$$\Delta p_x = 0$$

$$\Delta p_y = 0$$

$$\Delta p_z = 0,$$

or as a single vector equation,¹⁶

$$\Delta \mathbf{p} = 0.$$

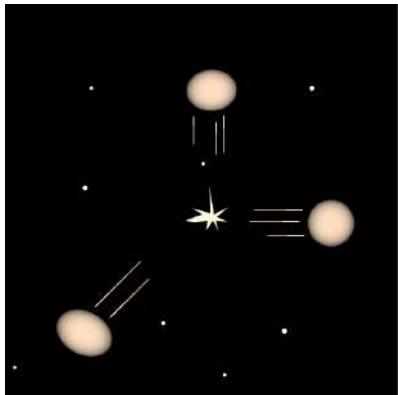
The following table summarizes some vector operations.

operation	definition
 vector 	$\sqrt{vector_x^2 + vector_y^2 + vector_z^2}$
vector + vector	Add component by component.
vector - vector	Subtract component by component.
vector · scalar	Multiply each component by the scalar.
vector / scalar	Divide each component by the scalar.

The first of these is called the *magnitude* of the vector; in one dimension, where a vector only has one component, it amounts to taking the absolute value, hence the similar notation.

self-check J

Translate the equations $F_x=ma_x$, $F_y=ma_y$, and $F_z=ma_z$ into a single equation in vector notation. ▷ Answer, p. 1060



h / Example 56.

An explosion

example 56

▷ Astronomers observe the planet Mars as the Martians fight a nuclear war. The Martian bombs are so powerful that they rip the planet into three separate pieces of liquefied rock, all having the same mass. If one fragment flies off with velocity components $v_{1x} = 0$, $v_{1y}=1.0 \times 10^4$ km/hr, and the second with $v_{2x}=1.0 \times 10^4$ km/hr, $v_{2y} = 0$, what is the magnitude of the third one's velocity?

▷ We work the problem in the center of mass frame, in which the planet initially had zero momentum. After the explosion, the vector sum of the momenta must still be zero. Vector addition can be done by adding components, so

$$mv_{1x} + mv_{2x} + mv_{3x} = 0$$

and

$$mv_{1y} + mv_{2y} + mv_{3y} = 0,$$

where we have used the same symbol *m* for all the terms, because the fragments all have the same mass. The masses can

¹⁶The zero here is really a zero *vector*, i.e., a vector whose components are all zero, so we should really represent it with a boldface 0. There's usually not much danger of confusion, however, so most books, including this one, don't use boldface for the zero vector.

be eliminated by dividing each equation by m , and we find

$$v_{3x} = -1.0 \times 10^4 \text{ km/hr}$$

$$v_{3y} = -1.0 \times 10^4 \text{ km/hr},$$

which gives a magnitude of

$$|\mathbf{v}_3| = \sqrt{v_{3x}^2 + v_{3y}^2}$$

$$= 1.4 \times 10^4 \text{ km/hr.}$$

A toppling box

example 57

If you place a box on a frictionless surface, it will fall over with a very complicated motion that is hard to predict in detail. We know, however, that its center of mass's motion is related to its momentum, and the rate at which momentum is transferred is the force. Moreover, we know that these relationships apply separately to each component. Let x and y be horizontal, and z vertical. There are two forces on the box, an upward force from the table and a downward gravitational force. Since both of these are along the z axis, p_z is the only component of the box's momentum that can change. We conclude that the center of mass travels vertically. This is true even if the box bounces and tumbles. [Based on an example by Kleppner and Kolenkow.]

Geometric representation of vectors

A vector in two dimensions can be easily visualized by drawing an arrow whose length represents its magnitude and whose direction represents its direction. The x component of a vector can then be visualized, j , as the length of the shadow it would cast in a beam of light projected onto the x axis, and similarly for the y component. Shadows with arrowheads pointing back against the direction of the positive axis correspond to negative components.

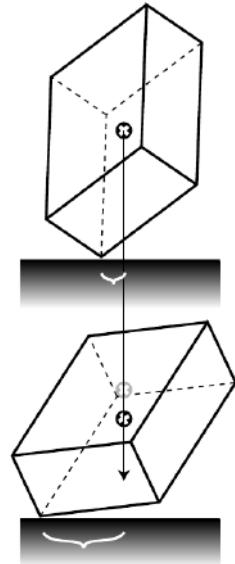
In this type of diagram, the negative of a vector is the vector with the same magnitude but in the opposite direction. Multiplying a vector by a scalar is represented by lengthening the arrow by that factor, and similarly for division.

self-check K

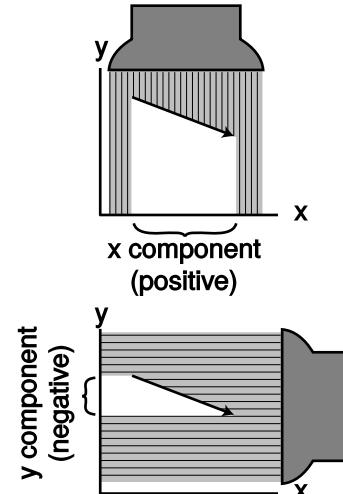
Given vector \mathbf{Q} represented by an arrow below, draw arrows representing the vectors $1.5\mathbf{Q}$ and $-\mathbf{Q}$.



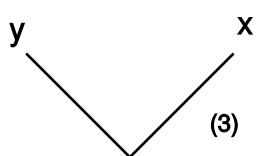
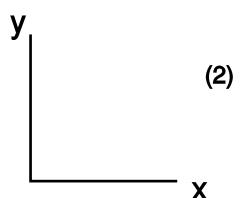
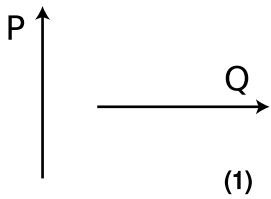
▷ Answer, p. 1060



i / Example 57.



j / The geometric interpretation of a vector's components.



k / Two vectors, 1, to which we apply the same operation in two different frames of reference, 2 and 3.

A useless vector operation

example 58

The way I've defined the various vector operations above aren't as arbitrary as they seem. There are many different vector operations that we could define, but only some of the possible definitions are mathematically useful. Consider the operation of multiplying two vectors component by component to produce a third vector:

$$\begin{aligned}R_x &= P_x Q_x \\R_y &= P_y Q_y \\R_z &= P_z Q_z\end{aligned}$$

As a simple example, we choose vectors **P** and **Q** to have length 1, and make them perpendicular to each other, as shown in figure k/1. If we compute the result of our new vector operation using the coordinate system shown in k/2, we find:

$$\begin{aligned}R_x &= 0 \\R_y &= 0 \\R_z &= 0\end{aligned}$$

The *x* component is zero because $P_x = 0$, the *y* component is zero because $Q_y = 0$, and the *z* component is of course zero because both vectors are in the *x*-*y* plane. However, if we carry out the same operations in coordinate system k/3, rotated 45 degrees with respect to the previous one, we find

$$\begin{aligned}R_x &= -1/2 \\R_y &= 1/2 \\R_z &= 0\end{aligned}$$

The operation's result depends on what coordinate system we use, and since the two versions of **R** have different lengths (one being zero and the other nonzero), they don't just represent the same answer expressed in two different coordinate systems. Such an operation will never be useful in physics, because experiments show physics works the same regardless of which way we orient the laboratory building! The useful vector operations, such as addition and scalar multiplication, are rotationally invariant, i.e., come out the same regardless of the orientation of the coordinate system.

All the vector techniques can be applied to any kind of vector, but the graphical representation of vectors as arrows is particularly natural for vectors that represent lengths and distances. We define a vector called **r** whose components are the coordinates of a particular point in space, *x*, *y*, and *z*. The $\Delta\mathbf{r}$ vector, whose components are Δx , Δy , and Δz , can then be used to represent motion that starts at

one point and ends at another. Adding two $\Delta\mathbf{r}$ vectors is interpreted as a trip with two legs: by computing the $\Delta\mathbf{r}$ vector going from point A to point B plus the vector from B to C, we find the vector that would have taken us directly from A to C.

Calculations with magnitude and direction

If you ask someone where Las Vegas is compared to Los Angeles, she is unlikely to say that the Δx is 290 km and the Δy is 230 km, in a coordinate system where the positive x axis is east and the y axis points north. She will probably say instead that it's 370 km to the northeast. If she was being precise, she might specify the direction as 38° counterclockwise from east. In two dimensions, we can always specify a vector's direction like this, using a single angle. A magnitude plus an angle suffice to specify everything about the vector. The following two examples show how we use trigonometry and the Pythagorean theorem to go back and forth between the x - y and magnitude-angle descriptions of vectors.

Finding magnitude and angle from components example 59

▷ Given that the $\Delta\mathbf{r}$ vector from LA to Las Vegas has $\Delta x=290$ km and $\Delta y=230$ km, how would we find the magnitude and direction of $\Delta\mathbf{r}$?

▷ We find the magnitude of $\Delta\mathbf{r}$ from the Pythagorean theorem:

$$\begin{aligned} |\Delta\mathbf{r}| &= \sqrt{\Delta x^2 + \Delta y^2} \\ &= 370 \text{ km} \end{aligned}$$

We know all three sides of the triangle, so the angle θ can be found using any of the inverse trig functions. For example, we know the opposite and adjacent sides, so

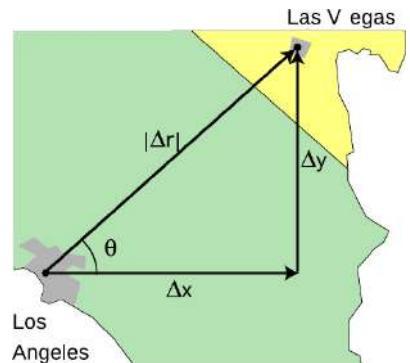
$$\begin{aligned} \theta &= \tan^{-1} \frac{\Delta y}{\Delta x} \\ &= 38^\circ. \end{aligned}$$

Finding the components from the magnitude and angle example 60

▷ Given that the straight-line distance from Los Angeles to Las Vegas is 370 km, and that the angle θ in the figure is 38° , how can the x and y components of the $\Delta\mathbf{r}$ vector be found?

▷ The sine and cosine of θ relate the given information to the information we wish to find:

$$\begin{aligned} \cos \theta &= \frac{\Delta x}{|\Delta\mathbf{r}|} \\ \sin \theta &= \frac{\Delta y}{|\Delta\mathbf{r}|} \end{aligned}$$



I / Example 59.

Solving for the unknowns gives

$$\begin{aligned}\Delta x &= |\Delta r| \cos \theta \\ &= 290 \text{ km} \\ \Delta y &= |\Delta r| \sin \theta \\ &= 230 \text{ km}\end{aligned}$$

The following example shows the correct handling of the plus and minus signs, which is usually the main cause of mistakes by students.

Negative components

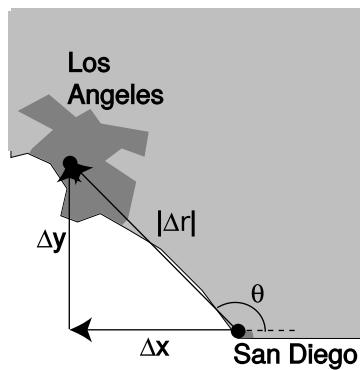
example 61

► San Diego is 120 km east and 150 km south of Los Angeles. An airplane pilot is setting course from San Diego to Los Angeles. At what angle should she set her course, measured counterclockwise from east, as shown in the figure?

► If we make the traditional choice of coordinate axes, with x pointing to the right and y pointing up on the map, then her Δx is negative, because her final x value is less than her initial x value. Her Δy is positive, so we have

$$\begin{aligned}\Delta x &= -120 \text{ km} \\ \Delta y &= 150 \text{ km.}\end{aligned}$$

If we work by analogy with the example 59, we get



m / Example 61.

$$\begin{aligned}\theta &= \tan^{-1} \frac{\Delta y}{\Delta x} \\ &= \tan^{-1} (-1.25) \\ &= -51^\circ.\end{aligned}$$

According to the usual way of defining angles in trigonometry, a negative result means an angle that lies clockwise from the x axis, which would have her heading for the Baja California. What went wrong? The answer is that when you ask your calculator to take the arctangent of a number, there are always two valid possibilities differing by 180° . That is, there are two possible angles whose tangents equal -1.25 :

$$\begin{aligned}\tan 129^\circ &= -1.25 \\ \tan (-51^\circ) &= -1.25\end{aligned}$$

Your calculator doesn't know which is the correct one, so it just picks one. In this case, the one it picked was the wrong one, and it was up to you to add 180° to it to find the right answer.

A shortcut

example 62

▷ A split second after nine o'clock, the hour hand on a clock dial has moved clockwise past the nine-o'clock position by some imperceptibly small angle ϕ . Let positive x be to the right and positive y up. If the hand, with length ℓ , is represented by a $\Delta\mathbf{r}$ vector going from the dial's center to the tip of the hand, find this vector's Δx .

▷ The following shortcut is the easiest way to work out examples like these, in which a vector's direction is known relative to one of the axes. We can tell that $\Delta\mathbf{r}$ will have a large, negative x component and a small, positive y . Since $\Delta x < 0$, there are really only two logical possibilities: either $\Delta x = -\ell \cos \phi$, or $\Delta x = -\ell \sin \phi$. Because ϕ is small, $\cos \phi$ is large and $\sin \phi$ is small. We conclude that $\Delta x = -\ell \cos \phi$.

A typical application of this technique to force vectors is given in example 71 on p. 209.

Addition of vectors given their components

The easiest type of vector addition is when you are in possession of the components, and want to find the components of their sum.

San Diego to Las Vegas

example 63

▷ Given the Δx and Δy values from the previous examples, find the Δx and Δy from San Diego to Las Vegas.

▷

$$\begin{aligned}\Delta x_{total} &= \Delta x_1 + \Delta x_2 \\ &= -120 \text{ km} + 290 \text{ km} \\ &= 170 \text{ km} \\ \Delta y_{total} &= \Delta y_1 + \Delta y_2 \\ &= 150 \text{ km} + 230 \text{ km} \\ &= 380\end{aligned}$$



n / Example 62.



o / Example 63.

Addition of vectors given their magnitudes and directions

In this case, you must first translate the magnitudes and directions into components, and then add the components.

Graphical addition of vectors

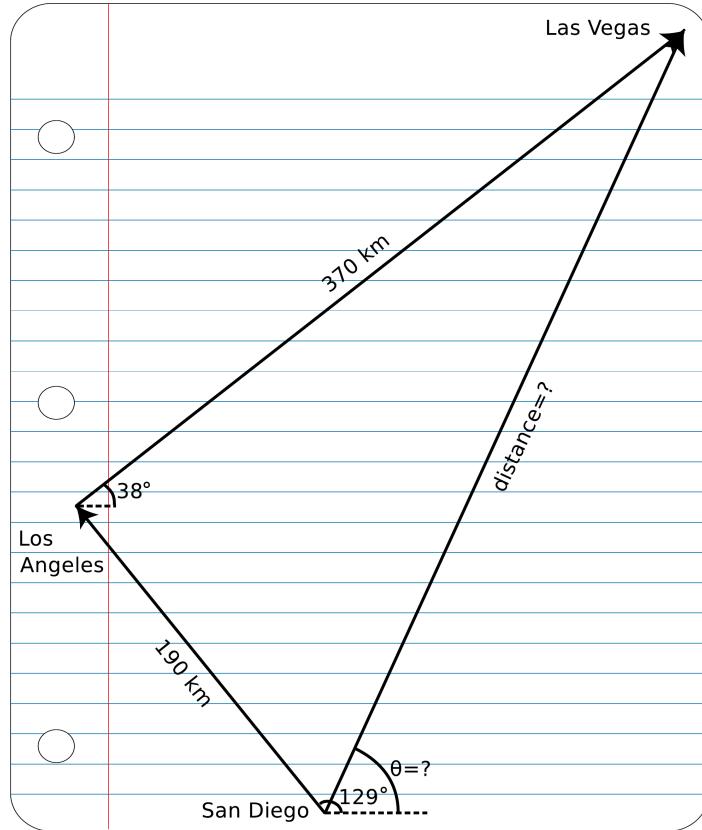
Often the easiest way to add vectors is by making a scale drawing on a piece of paper. This is known as graphical addition, as opposed to the analytic techniques discussed previously.

From San Diego to Las Vegas, graphically

example 64

▷ Given the magnitudes and angles of the $\Delta\mathbf{r}$ vectors from San

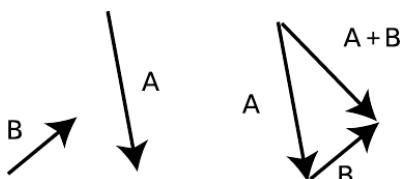
p / Example 64.



Diego to Los Angeles and from Los Angeles to Las Vegas, find the magnitude and angle of the $\Delta\mathbf{r}$ vector from San Diego to Las Vegas.

▷ Using a protractor and a ruler, we make a careful scale drawing, as shown in figure p on page 204. A scale of 1 cm \leftrightarrow 10 km was chosen for this solution. With a ruler, we measure the distance from San Diego to Las Vegas to be 3.8 cm, which corresponds to 380 km. With a protractor, we measure the angle θ to be 71° .

Even when we don't intend to do an actual graphical calculation with a ruler and protractor, it can be convenient to diagram the addition of vectors in this way, as shown in figure q. With $\Delta\mathbf{r}$ vectors, it intuitively makes sense to lay the vectors tip-to-tail and draw the sum vector from the tail of the first vector to the tip of the second vector. We can do the same when adding other vectors such as force vectors.



q / Adding vectors graphically by placing them tip-to-tail, like a train.

Unit vector notation

When we want to specify a vector by its components, it can be cumbersome to have to write the algebra symbol for each component:

$$\Delta x = 290 \text{ km}, \quad \Delta y = 230 \text{ km}$$

A more compact notation is to write

$$\Delta \mathbf{r} = (290 \text{ km})\hat{\mathbf{x}} + (230 \text{ km})\hat{\mathbf{y}},$$

where the vectors $\hat{\mathbf{x}}$, $\hat{\mathbf{y}}$, and $\hat{\mathbf{z}}$, called the unit vectors, are defined as the vectors that have magnitude equal to 1 and directions lying along the x , y , and z axes. In speech, they are referred to as “x-hat,” “y-hat,” and “z-hat.”

A slightly different, and harder to remember, version of this notation is unfortunately more prevalent. In this version, the unit vectors are called $\hat{\mathbf{i}}$, $\hat{\mathbf{j}}$, and $\hat{\mathbf{k}}$:

$$\Delta \mathbf{r} = (290 \text{ km})\hat{\mathbf{i}} + (230 \text{ km})\hat{\mathbf{j}}.$$

Applications to relative motion, momentum, and force

Vector addition is the correct way to generalize the one-dimensional concept of adding velocities in relative motion, as shown in the following example:

Velocity vectors in relative motion

example 65

▷ You wish to cross a river and arrive at a dock that is directly across from you, but the river's current will tend to carry you downstream. To compensate, you must steer the boat at an angle. Find the angle θ , given the magnitude, $|\mathbf{v}_{WL}|$, of the water's velocity relative to the land, and the maximum speed, $|\mathbf{v}_{BW}|$, of which the boat is capable relative to the water.

▷ The boat's velocity relative to the land equals the vector sum of its velocity with respect to the water and the water's velocity with respect to the land,

$$\mathbf{v}_{BL} = \mathbf{v}_{BW} + \mathbf{v}_{WL}.$$

If the boat is to travel straight across the river, i.e., along the y axis, then we need to have $v_{BL,x} = 0$. This x component equals the sum of the x components of the other two vectors,

$$v_{BL,x} = v_{BW,x} + v_{WL,x},$$

or

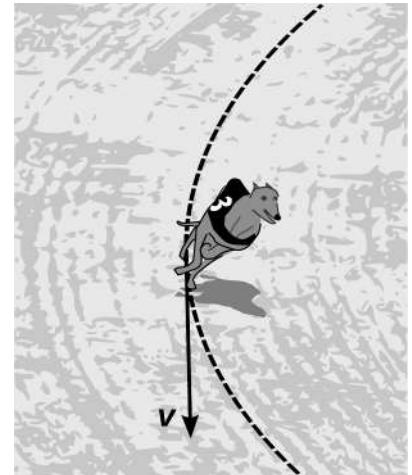
$$0 = -|v_{BW}| \sin \theta + |v_{WL}|.$$

Solving for θ , we find

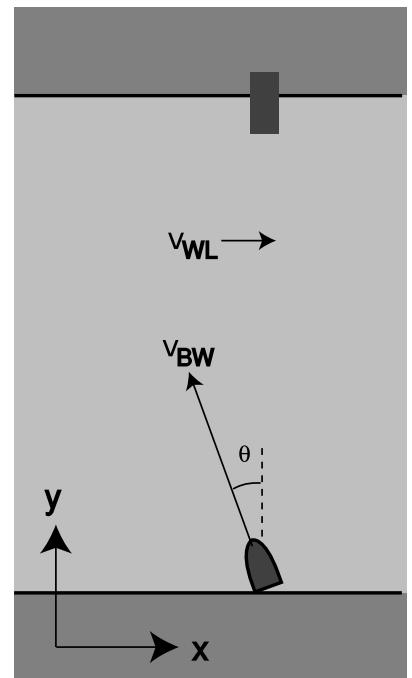
$$\sin \theta = |\mathbf{v}_{WL}| / |\mathbf{v}_{BW}|,$$

so

$$\theta = \sin^{-1} \frac{|\mathbf{v}_{WL}|}{|\mathbf{v}_{BW}|}.$$



r / The racing greyhound's velocity vector is in the direction of its motion, i.e., tangent to its curved path.



s / Example 65

How to generalize one-dimensional equations

example 66

- ▷ How can the one-dimensional relationships

$$p_{total} = m_{total} v_{cm}$$

and

$$x_{cm} = \frac{\sum_j m_j x_j}{\sum_j m_j}$$

be generalized to three dimensions?

- ▷ Momentum and velocity are vectors, since they have directions in space. Mass is a scalar. If we rewrite the first equation to show the appropriate quantities notated as vectors,

$$\mathbf{p}_{total} = m_{total} \mathbf{v}_{cm},$$

we get a valid mathematical operation, the multiplication of a vector by a scalar. Similarly, the second equation becomes

$$\mathbf{r}_{cm} = \frac{\sum_j m_j \mathbf{r}_j}{\sum_j m_j},$$

which is also valid. Each term in the sum on top contains a vector multiplied by a scalar, which gives a vector. Adding up all these vectors gives a vector, and dividing by the scalar sum on the bottom gives another vector.

This kind of wave-the-magic-wand-and-write-it-all-in-bold-face technique will always give the right generalization from one dimension to three, provided that the result makes sense mathematically — if you find yourself doing something nonsensical, such as adding a scalar to a vector, then you haven't found the generalization correctly.

Colliding coins

example 67

- ▷ Take two identical coins, put one down on a piece of paper, and slide the other across the paper, shooting it fairly rapidly so that it hits the target coin off-center. If you trace the initial and final positions of the coins, you can determine the directions of their momentum vectors after the collision. The angle between these vectors is always fairly close to, but a little less than, 90 degrees. Why is this?

- ▷ Let the velocity vector of the incoming coin be \mathbf{a} , and let the two outgoing velocity vectors be \mathbf{b} and \mathbf{c} . Since the masses are the same, conservation of momentum amounts to $\mathbf{a} = \mathbf{b} + \mathbf{c}$, which means that it has to be possible to assemble the three vectors into a triangle. If we assume that no energy is converted into heat and sound, then conservation of energy gives (discarding the common factor of $m/2$) $a^2 = b^2 + c^2$ for the magnitudes of the

three vectors. This is the Pythagorean theorem, which will hold only if the three vectors form a right triangle.

The fact that we observe the angle to be somewhat less than 90 degrees shows that the assumption used in the proof is only approximately valid: a little energy *is* converted into heat and sound. The opposite case would be a collision between two blobs of putty, where the maximum possible amount of energy is converted into heat and sound, the two blobs fly off together, giving an angle of zero between their momentum vectors. The real-life experiment interpolates between the ideal extremes of 0 and 90 degrees, but comes much closer to 90.

Force is a vector, and we add force vectors when more than one force acts on the same object.

Pushing a block up a ramp

example 68

▷ Figure t/1 shows a block being pushed up a frictionless ramp at constant speed by an applied force F_a . How much force is required, in terms of the block's mass, m , and the angle of the ramp, θ ?

▷ We analyzed this simple machine in example 38 on page 172 using the concept of work. Here we'll do it using vector addition of forces. Figure t/2 shows the other two forces acting on the block: a normal force, F_n , created by the ramp, and the gravitational force, F_g . Because the block is being pushed up at constant speed, it has zero acceleration, and the total force on it must be zero. In figure t/3, we position all the force vectors tip-to-tail for addition. Since they have to add up to zero, they must join up without leaving a gap, so they form a triangle. Using trigonometry we find

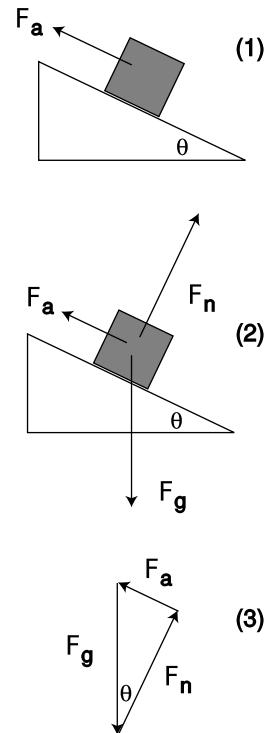
$$F_a = F_g \sin \theta \\ = mg \sin \theta.$$

Buoyancy, again

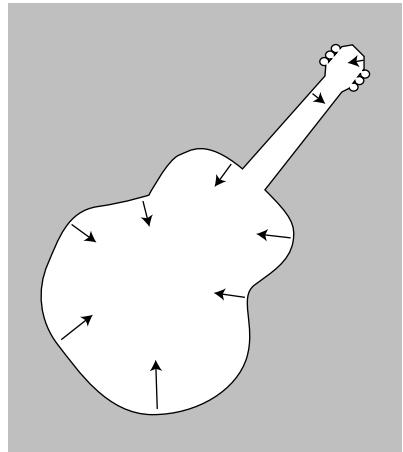
example 69

In example 10 on page 85, we found that the energy required to raise a cube immersed in a fluid is as if the cube's mass had been reduced by an amount equal to the mass of the fluid that otherwise would have been in the volume it occupies (Archimedes' principle). From the energy perspective, this effect occurs because raising the cube allows a certain amount of fluid to move downward, and the decreased gravitational energy of the fluid tends to offset the increased gravitational energy of the cube. The proof given there, however, could not easily be extended to other shapes.

Thinking in terms of force rather than energy, it becomes easier to give a proof that works for any shape. A certain upward force is



t / Example 68.



u / Archimedes' principle works regardless of whether the object is a cube. The fluid makes a force on every square millimeter of the object's surface.

needed to support the object in figure u. If this force was applied, then the object would be in equilibrium: the vector sum of all the forces acting on it would be zero. These forces are \mathbf{F}_a , the upward force just mentioned, \mathbf{F}_g , the downward force of gravity, and \mathbf{F}_f , the total force from the fluid:

$$\mathbf{F}_a + \mathbf{F}_g + \mathbf{F}_f = 0$$

Since the fluid is under more pressure at a greater depth, the part of the fluid underneath the object tends to make more force than the part above, so the fluid tends to help support the object.

Now suppose the object was removed, and instantly replaced with an equal volume of fluid. The new fluid would be in equilibrium without any force applied to hold it up, so

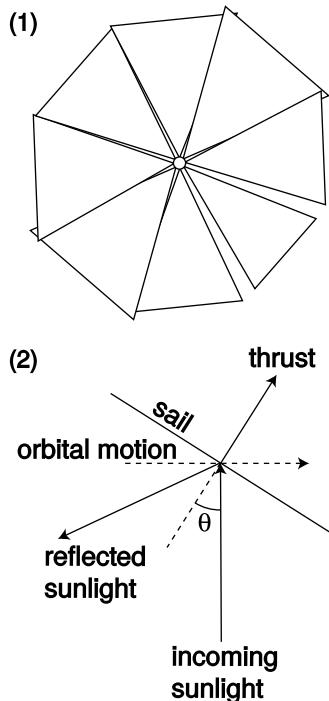
$$\mathbf{F}_{gf} + \mathbf{F}_f = 0,$$

where \mathbf{F}_{gf} , the weight of the fluid, is not the same as \mathbf{F}_g , the weight of the object, but \mathbf{F}_f is the same as before, since the pressure of the surrounding fluid is the same as before at any particular depth. We therefore have

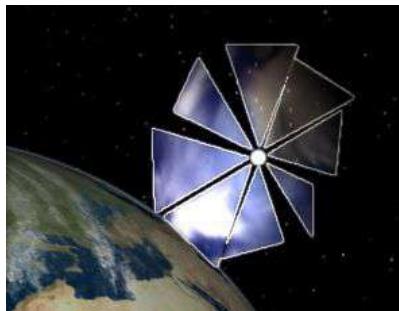
$$\mathbf{F}_a = -(\mathbf{F}_g - \mathbf{F}_{gf}),$$

which is Archimedes' principle in terms of force: the force required to support the object is lessened by an amount equal to the weight of the fluid that would have occupied its volume.

By the way, the word "pressure" that I threw around casually in the preceding example has a precise technical definition: force per unit area. The SI units of pressure are N/m^2 , which can be abbreviated as pascals, $1 \text{ Pa} = 1 \text{ N/m}^2$. Atmospheric pressure is about 100 kPa. By applying the equation $\mathbf{F}_g + \mathbf{F}_f = 0$ to the top and bottom surfaces of a cubical volume of fluid, one can easily prove that the difference in pressure between two different depths is $\Delta P = \rho g \Delta y$. (In physics, "fluid" can refer to either a gas or a liquid.) Pressure is discussed in more detail in chapter 5.



v / Example 70.



w / An artist's rendering of what Cosmos 1 would have looked like in orbit.

A solar sail

example 70

A solar sail, figure v/1, allows a spacecraft to get its thrust without using internal stores of energy or having to carry along mass that it can shove out the back like a rocket. Sunlight strikes the sail and bounces off, transferring momentum to the sail. A working 30-meter-diameter solar sail, *Cosmos 1*, was built by an American company, and was supposed to be launched into orbit aboard a Russian booster launched from a submarine, but launch attempts in 2001 and 2005 both failed.

In this example, we will calculate the optimal orientation of the sail, assuming that "optimal" means changing the vehicle's energy as rapidly as possible. For simplicity, we model the complicated shape of the sail's surface as a disk, seen edge-on in figure v/2, and we assume that the craft is in a nearly circular orbit

around the sun, hence the 90-degree angle between the direction of motion and the incoming sunlight. We assume that the sail is 100% reflective. The orientation of the sail is specified using the angle θ between the incoming rays of sunlight and the perpendicular to the sail. In other words, $\theta=0$ if the sail is catching the sunlight full-on, while $\theta=90^\circ$ means that the sail is edge-on to the sun.

Conservation of momentum gives

$$\mathbf{p}_{light,i} = \mathbf{p}_{light,f} + \Delta\mathbf{p}_{sail},$$

where $\Delta\mathbf{p}_{sail}$ is the change in momentum picked up by the sail. Breaking this down into components, we have

$$0 = p_{light,f,x} + \Delta p_{sail,x} \quad \text{and}$$

$$p_{light,i,y} = p_{light,f,y} + \Delta p_{sail,y}.$$

As in example 53 on page 193, the component of the force that is directly away from the sun (up in figure v/2) doesn't change the energy of the craft, so we only care about $\Delta p_{sail,x}$, which equals $-p_{light,f,x}$. The outgoing light ray forms an angle of 2θ with the negative y axis, or $270^\circ - 2\theta$ measured counterclockwise from the x axis, so the useful thrust depends on $-\cos(270^\circ - 2\theta) = \sin 2\theta$.

However, this is all assuming a given amount of light strikes the sail. During a certain time period, the amount of sunlight striking the sail depends on the cross-sectional area the sail presents to the sun, which is proportional to $\cos \theta$. For $\theta=90^\circ$, $\cos \theta$ equals zero, since the sail is edge-on to the sun.

Putting together these two factors, the useful thrust is proportional to $\sin 2\theta \cos \theta$, and this quantity is maximized for $\theta \approx 35^\circ$. A counterintuitive fact about this maneuver is that as the spacecraft spirals outward, its total energy (kinetic plus gravitational) increases, but its kinetic energy actually decreases!

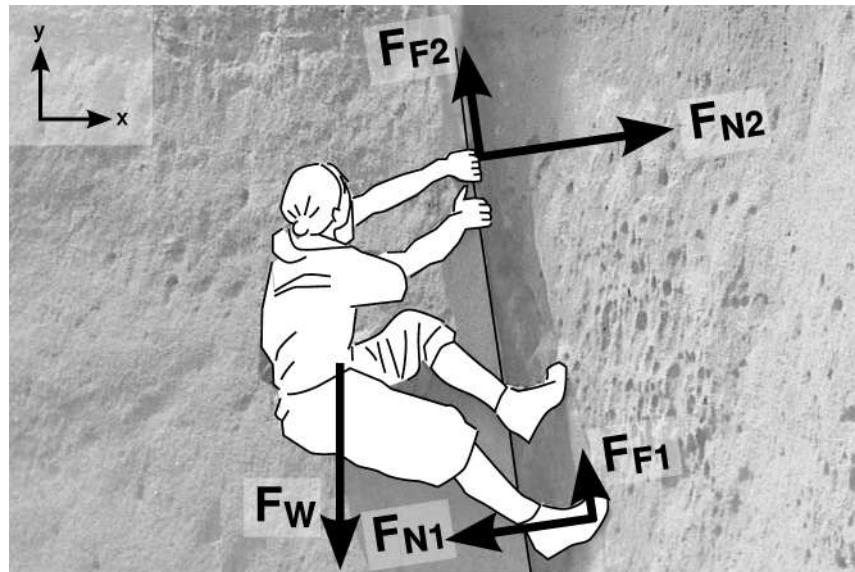
A layback

example 71

The figure shows a rock climber using a technique called a layback. He can make the normal forces \mathbf{F}_{N1} and \mathbf{F}_{N2} large, which has the side-effect of increasing the frictional forces \mathbf{F}_{F1} and \mathbf{F}_{F2} , so that he doesn't slip down due to the gravitational (weight) force \mathbf{F}_W . The purpose of the problem is not to analyze all of this in detail, but simply to practice finding the components of the forces based on their magnitudes. To keep the notation simple, let's write F_{N1} for $|\mathbf{F}_{N1}|$, etc. The crack overhangs by a small, positive angle $\theta \approx 9^\circ$.

In this example, we determine the x component of \mathbf{F}_{N1} . The other nine components are left as an exercise to the reader (problem 81, p. 239).

The easiest method is the one demonstrated in example 62 on p. 203. Casting vector \mathbf{F}_{N1} 's shadow on the ground, we can tell that it would point to the left, so its x component is negative. The only two possibilities for its x component are therefore $-F_{N1} \cos \theta$ or $-F_{N1} \sin \theta$. We expect this force to have a large x component and a much smaller y . Since θ is small, $\cos \theta \approx 1$, while $\sin \theta$ is small. Therefore the x component must be $-F_{N1} \cos \theta$.



x / Example 71 and problem 81 on p. 239.

Discussion Questions

A An object goes from one point in space to another. After it arrives at its destination, how does the magnitude of its $\Delta\mathbf{r}$ vector compare with the distance it traveled?

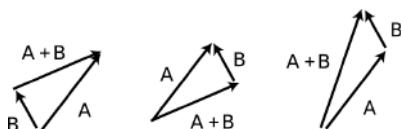
B In several examples, I've dealt with vectors having negative components. Does it make sense as well to talk about negative and positive vectors?

C If you're doing *graphical* addition of vectors, does it matter which vector you start with and which vector you start from the other vector's tip?

D If you add a vector with magnitude 1 to a vector of magnitude 2, what magnitudes are possible for the vector sum?

E Which of these examples of vector addition are correct, and which are incorrect?

F Is it possible for an airplane to maintain a constant velocity vector but not a constant $|\mathbf{v}|$? How about the opposite – a constant $|\mathbf{v}|$ but not a constant velocity vector? Explain.



y / Discussion question E.

G New York and Rome are at about the same latitude, so the earth's rotation carries them both around nearly the same circle. Do the two cities have the same velocity vector (relative to the center of the earth)? If not, is there any way for two cities to have the same velocity vector?

H The figure shows a roller coaster car rolling down and then up under the influence of gravity. Sketch the car's velocity vectors and acceleration vectors. Pick an interesting point in the motion and sketch a set of force vectors acting on the car whose vector sum could have resulted in the right acceleration vector.

I The following is a question commonly asked by students:

"Why does the force vector always have to point in the same direction as the acceleration vector? What if you suddenly decide to change your force on an object, so that your force is no longer pointing in the same direction that the object is accelerating?"

What misunderstanding is demonstrated by this question? Suppose, for example, a spacecraft is blasting its rear main engines while moving forward, then suddenly begins firing its sideways maneuvering rocket as well. What does the student think Newton's laws are predicting?

J Debug the following *incorrect* solutions to this vector addition problem.

Problem: Freddi Fish™ swims 5.0 km northeast, and then 12.0 km in the direction 55 degrees west of south. How far does she end up from her starting point, and in what direction is she from her starting point?

Incorrect solution #1:

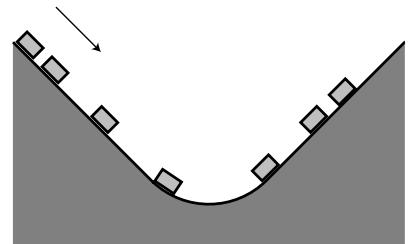
$$5.0 \text{ km} + 12.0 \text{ km} = 17.0 \text{ km}$$

Incorrect solution #2:

$$\sqrt{(5.0 \text{ km})^2 + (12.0 \text{ km})^2} = 13.0 \text{ km}$$

Incorrect solution #3:

Let \mathbf{A} and \mathbf{B} be her two $\Delta\mathbf{r}$ vectors, and let $\mathbf{C} = \mathbf{A} + \mathbf{B}$. Then



z / Discussion question H.

$$A_x = (5.0 \text{ km}) \cos 45^\circ = 3.5 \text{ km}$$

$$B_x = (12.0 \text{ km}) \cos 55^\circ = 6.9 \text{ km}$$

$$A_y = (5.0 \text{ km}) \sin 45^\circ = 3.5 \text{ km}$$

$$B_y = (12.0 \text{ km}) \sin 55^\circ = 9.8 \text{ km}$$

$$\begin{aligned} C_x &= A_x + B_x \\ &= 10.4 \text{ km} \end{aligned}$$

$$\begin{aligned} C_y &= A_y + B_y \\ &= 13.3 \text{ km} \end{aligned}$$

$$\begin{aligned} |\mathbf{C}| &= \sqrt{C_x^2 + C_y^2} \\ &= 16.9 \text{ km} \end{aligned}$$

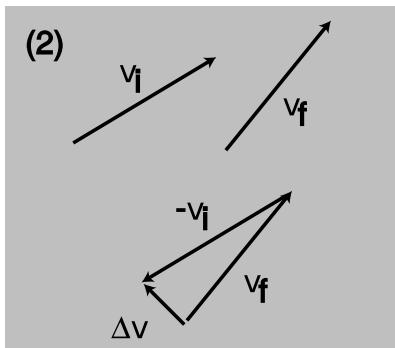
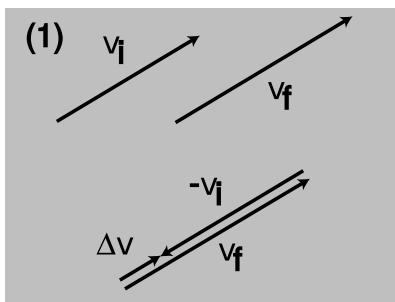
$$\begin{aligned} \text{direction} &= \tan^{-1}(13.3/10.4) \\ &= 52^\circ \text{ north of east} \end{aligned}$$

Incorrect solution #4:
(same notation as above)

$$\begin{aligned}
 A_x &= (5.0 \text{ km}) \cos 45^\circ = 3.5 \text{ km} \\
 B_x &= -(12.0 \text{ km}) \cos 55^\circ = -6.9 \text{ km} \\
 A_y &= (5.0 \text{ km}) \sin 45^\circ = 3.5 \text{ km} \\
 B_y &= -(12.0 \text{ km}) \sin 55^\circ = -9.8 \text{ km} \\
 C_x &= A_x + B_x \\
 &= -3.4 \text{ km} \\
 C_y &= A_y + B_y \\
 &= -6.3 \text{ km} \\
 |\mathbf{C}| &= \sqrt{C_x^2 + C_y^2} \\
 &= 7.2 \text{ km} \\
 \text{direction} &= \tan^{-1}(-6.3 / -3.4) \\
 &= 62^\circ \text{ north of east}
 \end{aligned}$$

Incorrect solution #5:
(same notation as above)

$$\begin{aligned}
 A_x &= (5.0 \text{ km}) \cos 45^\circ = 3.5 \text{ km} \\
 B_x &= -(12.0 \text{ km}) \sin 55^\circ = -9.8 \text{ km} \\
 A_y &= (5.0 \text{ km}) \sin 45^\circ = 3.5 \text{ km} \\
 B_y &= -(12.0 \text{ km}) \cos 55^\circ = -6.9 \text{ km} \\
 C_x &= A_x + B_x \\
 &= -6.3 \text{ km} \\
 C_y &= A_y + B_y \\
 &= -3.4 \text{ km} \\
 |\mathbf{C}| &= \sqrt{C_x^2 + C_y^2} \\
 &= 7.2 \text{ km} \\
 \text{direction} &= \tan^{-1}(-3.4 / -6.3) \\
 &= 28^\circ \text{ north of east}
 \end{aligned}$$



aa / Visualizing the acceleration vector.

3.4.4 Calculus with vectors

Differentiation

In one dimension, we define the velocity as the derivative of the position with respect to time, and we can think of the derivative as what we get when we calculate $\Delta x / \Delta t$ for very short time intervals. The quantity $\Delta x = x_f - x_i$ is calculated by subtraction. In three dimensions, x becomes \mathbf{r} , and the $\Delta \mathbf{r}$ vector is calculated by *vector* subtraction, $\Delta \mathbf{r} = \mathbf{r}_f - \mathbf{r}_i$. Vector subtraction is defined component by component, so when we take the derivative of a vector, this means we end up taking the derivative component by component,

$$v_x = \frac{dx}{dt}, \quad v_y = \frac{dy}{dt}, \quad v_z = \frac{dz}{dt}$$

or

$$\frac{d\mathbf{r}}{dt} = \frac{dx}{dt}\hat{\mathbf{x}} + \frac{dy}{dt}\hat{\mathbf{y}} + \frac{dz}{dt}\hat{\mathbf{z}}.$$

All of this reasoning applies equally well to any derivative of a vector, so for instance we can take the second derivative,

$$a_x = \frac{dv_x}{dt}, \quad a_y = \frac{dv_y}{dt}, \quad a_z = \frac{dv_z}{dt}$$

or

$$\frac{d\mathbf{v}}{dt} = \frac{dv_x}{dt}\hat{\mathbf{x}} + \frac{dv_y}{dt}\hat{\mathbf{y}} + \frac{dv_z}{dt}\hat{\mathbf{z}}.$$

A counterintuitive consequence of this is that the acceleration vector does not need to be in the same direction as the motion. The velocity vector points in the direction of motion, but by Newton's second law, $\mathbf{a} = \mathbf{F}/m$, the acceleration vector points in the same direction as the force, not the motion. This is easiest to understand if we take velocity vectors from two different moments in the motion, and visualize subtracting them graphically to make a $\Delta\mathbf{v}$ vector. The direction of the $\Delta\mathbf{v}$ vector tells us the direction of the acceleration vector as well, since the derivative $d\mathbf{v}/dt$ can be approximated as $\Delta\mathbf{v}/\Delta t$. As shown in figure aa/1, a change in the magnitude of the velocity vector implies an acceleration that is in the direction of motion. A change in the direction of the velocity vector produces an acceleration perpendicular to the motion, aa/2.

Circular motion

example 72

- ▷ An object moving in a circle of radius r in the x - y plane has

$$x = r \cos \omega t \quad \text{and} \\ y = r \sin \omega t,$$

where ω is the number of radians traveled per second, and the positive or negative sign indicates whether the motion is clockwise or counterclockwise. What is its acceleration?

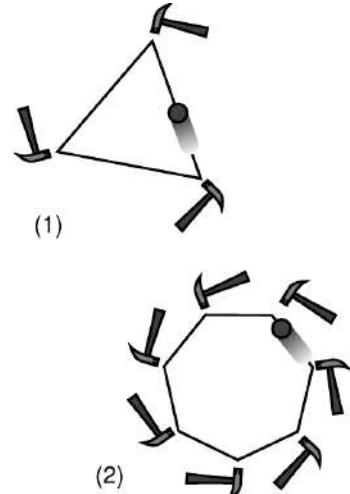
- ▷ The components of the velocity are

$$v_x = -\omega r \sin \omega t \quad \text{and} \\ v_y = \omega r \cos \omega t,$$

and for the acceleration we have

$$a_x = -\omega^2 r \cos \omega t \quad \text{and} \\ a_y = -\omega^2 r \sin \omega t.$$

The acceleration vector has cosines and sines in the same places as the \mathbf{v} vector, but with minus signs in front, so it points in the opposite direction, i.e., toward the center of the circle. By Newton's second law, $\mathbf{a}=\mathbf{F}/m$, this shows that the force must be inward as well; without this force, the object would fly off straight.



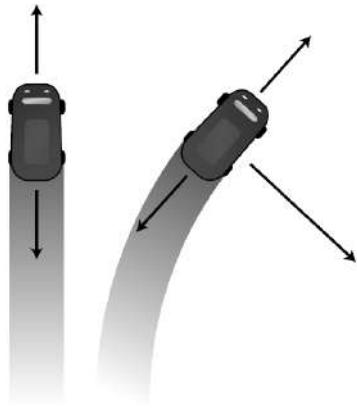
ab / This figure shows an intuitive justification for the fact proved mathematically in the example, that the direction of the force and acceleration in circular motion is inward. The heptagon, 2, is a better approximation to a circle than the triangle, 1. To make an infinitely good approximation to circular motion, we would need to use an infinitely large number of infinitesimal taps, which would amount to a steady inward force.

The magnitude of the acceleration is

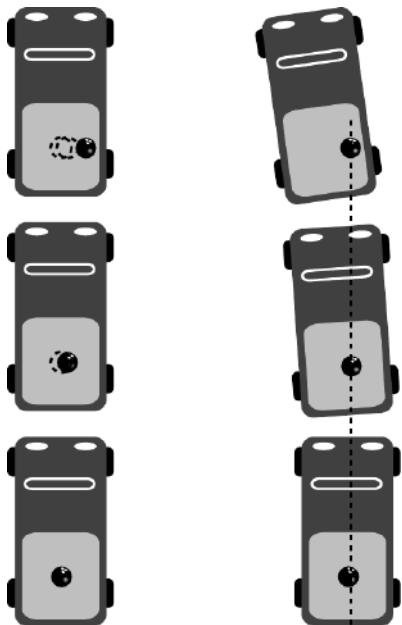
$$|\mathbf{a}| = \sqrt{a_x^2 + a_y^2} \\ = \omega^2 r.$$

It makes sense that ω is squared, since reversing the sign of ω corresponds to reversing the direction of motion, but the acceleration is toward the center of the circle, regardless of whether the motion is clockwise or counterclockwise. This result can also be rewritten in the form

$$|\mathbf{a}| = \frac{|\mathbf{v}|^2}{r}.$$



ac / The total force in the forward-backward direction is zero in both cases.



ad / There is no outward force on the bowling ball, but in the noninertial frame it seems like one exists.

Although I've relegated the results $a = \omega^2 r = |\mathbf{v}|^2/r$ to an example because they are a straightforward corollary of more general principles already developed, they are important and useful enough to record for later use. These results are counterintuitive as well. Until Newton, physicists and laypeople alike had assumed that the planets would need a force to push them *forward* in their orbits. Figure ab may help to make it more plausible that only an inward force is required. A forward force might be needed in order to cancel out a backward force such as friction, ac, but the total force in the forward-backward direction needs to be exactly zero for constant-speed motion. When you are in a car undergoing circular motion, there is also a strong illusion of an *outward* force. But what object could be making such a force? The car's seat makes an inward force on you, not an outward one. There is no object that could be exerting an outward force on your body. In reality, this force is an illusion that comes from our brain's intuitive efforts to interpret the situation within a noninertial frame of reference. As shown in figure ad, we can describe everything perfectly well in an inertial frame of reference, such as the frame attached to the sidewalk. In such a frame, the bowling ball goes straight because there is *no* force on it. The wall of the truck's bed hits the ball, not the other way around.

Integration

An integral is really just a sum of many infinitesimally small terms. Since vector addition is defined in terms of addition of the components, an integral of a vector quantity is found by doing integrals component by component.

Projectile motion

example 73

- ▷ Find the motion of an object whose acceleration vector is constant, for instance a projectile moving under the influence of gravity.
- ▷ We integrate the acceleration to get the velocity, and then integrate the velocity to get the position as a function of time. Doing

this to the x component of the acceleration, we find

$$\begin{aligned}x &= \int \left(\int a_x dt \right) dt \\&= \int (a_x t + v_{x0}) dt,\end{aligned}$$

where v_{x0} is a constant of integration, and

$$x = \frac{1}{2} a_x t^2 + v_{x0} t + x_0.$$

Similarly, $y = (1/2)a_y t^2 + v_{y0} t + y_0$ and $z = (1/2)a_z t^2 + v_{z0} t + z_0$. Once one has gained a little confidence, it becomes natural to do the whole thing as a single vector integral,

$$\begin{aligned}\mathbf{r} &= \int \left(\int \mathbf{a} dt \right) dt \\&= \int (\mathbf{a}t + \mathbf{v}_0) dt \\&= \frac{1}{2} \mathbf{a}t^2 + \mathbf{v}_0 t + \mathbf{r}_0,\end{aligned}$$

where now the constants of integration are vectors.

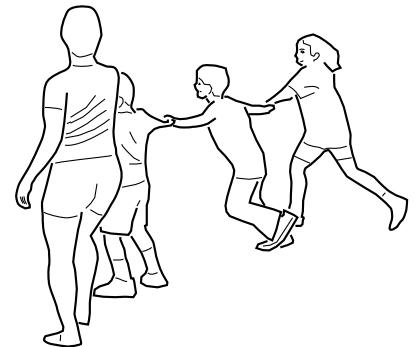
Discussion Questions

A In the game of crack the whip, a line of people stand holding hands, and then they start sweeping out a circle. One person is at the center, and rotates without changing location. At the opposite end is the person who is running the fastest, in a wide circle. In this game, someone always ends up losing their grip and flying off. Suppose the person on the end loses her grip. What path does she follow as she goes flying off? Draw an overhead view. (Assume she is going so fast that she is really just trying to put one foot in front of the other fast enough to keep from falling; she is not able to get any significant horizontal force between her feet and the ground.)

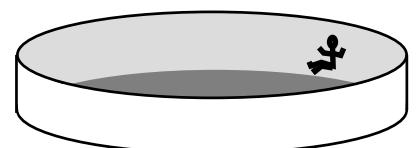
B Suppose the person on the outside is still holding on, but feels that she may lose her grip at any moment. What force or forces are acting on her, and in what directions are they? (We are not interested in the vertical forces, which are the earth's gravitational force pulling down, and the ground's normal force pushing up.) Make a table in the format shown in subsection 3.2.6.

C Suppose the person on the outside is still holding on, but feels that she may lose her grip at any moment. What is wrong with the following analysis of the situation? "The person whose hand she's holding exerts an inward force on her, and because of Newton's third law, there's an equal and opposite force acting outward. That outward force is the one she feels throwing her outward, and the outward force is what might make her go flying off, if it's strong enough."

D If the only force felt by the person on the outside is an inward force, why doesn't she go straight in?



ae / Discussion question A.



af / Discussion question E.

E In the amusement park ride shown in the figure, the cylinder spins faster and faster until the customer can pick her feet up off the floor without falling. In the old Coney Island version of the ride, the floor actually dropped out like a trap door, showing the ocean below. (There is also a version in which the whole thing tilts up diagonally, but we're discussing the version that stays flat.) If there is no outward force acting on her, why does she stick to the wall? Analyze all the forces on her.

F What is an example of circular motion where the inward force is a normal force? What is an example of circular motion where the inward force is friction? What is an example of circular motion where the inward force is the sum of more than one force?

G Does the acceleration vector always change continuously in circular motion? The velocity vector?

H A certain amount of force is needed to provide the acceleration of circular motion. What if we are exerting a force perpendicular to the direction of motion in an attempt to make an object trace a circle of radius r , but the force isn't as big as $m|\mathbf{v}|^2/r$?

I Suppose a rotating space station is built that gives its occupants the illusion of ordinary gravity. What happens when a person in the station lets go of a ball? What happens when she throws a ball straight "up" in the air (i.e., towards the center)?

3.4.5 The dot product

How would we generalize the mechanical work equation $dE = F dx$ to three dimensions? Energy is a scalar, but force and distance are vectors, so it might seem at first that the kind of "magic-wand" generalization discussed on page 206 failed here, since we don't know of any way to multiply two vectors together to get a scalar. Actually, this is Nature giving us a hint that there is such a multiplication operation waiting for us to invent it, and since Nature is simple, we can be assured that this operation will work just fine in any situation where a similar generalization is required.

How should this operation be defined? Let's consider what we would get by performing this operation on various combinations of the unit vectors $\hat{\mathbf{x}}$, $\hat{\mathbf{y}}$, and $\hat{\mathbf{z}}$. The conventional notation for the operation is to put a dot, \cdot , between the two vectors, and the operation is therefore called the *dot product*. Rotational invariance requires that we handle the three coordinate axes in the same way, without giving special treatment to any of them, so we must have $\hat{\mathbf{x}} \cdot \hat{\mathbf{x}} = \hat{\mathbf{y}} \cdot \hat{\mathbf{y}} = \hat{\mathbf{z}} \cdot \hat{\mathbf{z}}$ and $\hat{\mathbf{x}} \cdot \hat{\mathbf{y}} = \hat{\mathbf{y}} \cdot \hat{\mathbf{z}} = \hat{\mathbf{z}} \cdot \hat{\mathbf{x}}$. This is supposed to be a way of generalizing ordinary multiplication, so for consistency with the property $1 \times 1 = 1$ of ordinary numbers, the result of multiplying a magnitude-one vector by itself had better be the scalar 1, so $\hat{\mathbf{x}} \cdot \hat{\mathbf{x}} = \hat{\mathbf{y}} \cdot \hat{\mathbf{y}} = \hat{\mathbf{z}} \cdot \hat{\mathbf{z}} = 1$. Furthermore, there is no way to satisfy rotational invariance unless we define the mixed products to be zero, $\hat{\mathbf{x}} \cdot \hat{\mathbf{y}} = \hat{\mathbf{y}} \cdot \hat{\mathbf{z}} = \hat{\mathbf{z}} \cdot \hat{\mathbf{x}} = 0$; for example, a 90-degree rotation of our frame of reference about the z axis reverses the sign of $\hat{\mathbf{x}} \cdot \hat{\mathbf{y}}$, but rotational invariance requires that $\hat{\mathbf{x}} \cdot \hat{\mathbf{y}}$ produce the

same result either way, and zero is the only number that stays the same when we reverse its sign. Establishing these six products of unit vectors suffices to define the operation in general, since any two vectors that we want to multiply can be broken down into components, e.g., $(2\hat{\mathbf{x}} + 3\hat{\mathbf{z}}) \cdot \hat{\mathbf{z}} = 2\hat{\mathbf{x}} \cdot \hat{\mathbf{z}} + 3\hat{\mathbf{z}} \cdot \hat{\mathbf{z}} = 0 + 3 = 3$. Thus by requiring rotational invariance and consistency with multiplication of ordinary numbers, we find that there is only one possible way to define a multiplication operation on two vectors that gives a scalar as the result.¹⁷ The dot product has all of the properties we normally associate with multiplication, except that there is no “dot division.”

Dot product in terms of components example 74

If we know the components of any two vectors \mathbf{b} and \mathbf{c} , we can find their dot product:

$$\begin{aligned}\mathbf{b} \cdot \mathbf{c} &= (b_x\hat{\mathbf{x}} + b_y\hat{\mathbf{y}} + b_z\hat{\mathbf{z}}) \cdot (c_x\hat{\mathbf{x}} + c_y\hat{\mathbf{y}} + c_z\hat{\mathbf{z}}) \\ &= b_x c_x + b_y c_y + b_z c_z.\end{aligned}$$

Magnitude expressed with a dot product example 75

If we take the dot product of any vector \mathbf{b} with itself, we find

$$\begin{aligned}\mathbf{b} \cdot \mathbf{b} &= (b_x\hat{\mathbf{x}} + b_y\hat{\mathbf{y}} + b_z\hat{\mathbf{z}}) \cdot (b_x\hat{\mathbf{x}} + b_y\hat{\mathbf{y}} + b_z\hat{\mathbf{z}}) \\ &= b_x^2 + b_y^2 + b_z^2,\end{aligned}$$

so its magnitude can be expressed as

$$|\mathbf{b}| = \sqrt{\mathbf{b} \cdot \mathbf{b}}.$$

We will often write b^2 to mean $\mathbf{b} \cdot \mathbf{b}$, when the context makes it clear what is intended. For example, we could express kinetic energy as $(1/2)m|\mathbf{v}|^2$, $(1/2)m\mathbf{v} \cdot \mathbf{v}$, or $(1/2)m\mathbf{v}^2$. In the third version, nothing but context tells us that \mathbf{v} really stands for the magnitude of some vector \mathbf{v} .

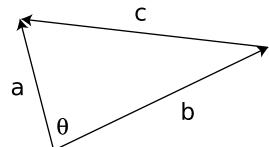
Geometric interpretation example 76

In figure ag, vectors \mathbf{a} , \mathbf{b} , and \mathbf{c} represent the sides of a triangle, and $\mathbf{a} = \mathbf{b} + \mathbf{c}$. The law of cosines gives

$$|\mathbf{c}|^2 = |\mathbf{a}|^2 + |\mathbf{b}|^2 - 2|\mathbf{a}||\mathbf{b}| \cos \theta.$$

Using the result of example 75, we can also write this as

$$\begin{aligned}|\mathbf{c}|^2 &= \mathbf{c} \cdot \mathbf{c} \\ &= (\mathbf{a} - \mathbf{b}) \cdot (\mathbf{a} - \mathbf{b}) \\ &= \mathbf{a} \cdot \mathbf{a} + \mathbf{b} \cdot \mathbf{b} - 2\mathbf{a} \cdot \mathbf{b}.\end{aligned}$$



ag / The geometric interpretation of the dot product.

¹⁷There is, however, a different operation, discussed in the next chapter, which multiplies two vectors to give a vector.

Matching up terms in these two expressions, we find

$$\mathbf{a} \cdot \mathbf{b} = |\mathbf{a}| |\mathbf{b}| \cos \theta,$$

which is a geometric interpretation for the dot product.

The result of example 76 is very useful. It gives us a way to find the angle between two vectors if we know their components. It can be used to show that the dot product of any two perpendicular vectors is zero. It also leads to a nifty proof that the dot product is rotationally invariant — up until now I've only proved that if a rotationally invariant product exists, the dot product is it — because angles and lengths aren't affected by a rotation, so the right side of the equation is rotationally invariant, and therefore so is the left side.

I introduced the whole discussion of the dot product by way of generalizing the equation $dE = F dx$ to three dimensions. In terms of a dot product, we have

$$dE = \mathbf{F} \cdot d\mathbf{r}.$$

If \mathbf{F} is a constant, integrating both sides gives

$$\Delta E = \mathbf{F} \cdot \Delta \mathbf{r}.$$

(If that step seemed like black magic, try writing it out in terms of components.) If the force is perpendicular to the motion, as in figure ah, then the work done is zero. The pack horse is doing work within its own body, but is not doing work on the pack.

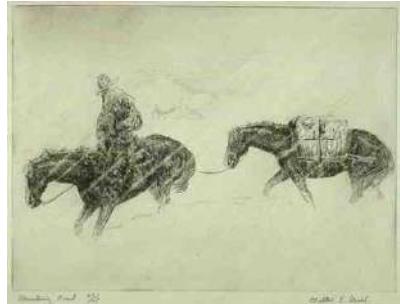
Pushing a lawnmower

example 77

▷ I push a lawnmower with a force $\mathbf{F} = (110 \text{ N})\hat{\mathbf{x}} - (40 \text{ N})\hat{\mathbf{y}}$, and the total distance I travel is $(100 \text{ m})\hat{\mathbf{x}}$. How much work do I do?

▷ The dot product is $11000 \text{ N}\cdot\text{m} = 11000 \text{ J}$.

A good application of the dot product is to allow us to write a simple, streamlined proof of separate conservation of the momentum components. (You can skip the proof without losing the continuity of the text.) The argument is a generalization of the one-dimensional proof on page 132, and makes the same assumption about the type of system of particles we're dealing with. The kinetic energy of one of the particles is $(1/2)m\mathbf{v} \cdot \mathbf{v}$, and when we transform into a different frame of reference moving with velocity \mathbf{u} relative to the original frame, the one-dimensional rule $v \rightarrow v + u$ turns into vector addition, $\mathbf{v} \rightarrow \mathbf{v} + \mathbf{u}$. In the new frame of reference, the kinetic energy is $(1/2)m(\mathbf{v} + \mathbf{u}) \cdot (\mathbf{v} + \mathbf{u})$. For a system of n particles, we



ah / Breaking trail, by Walter E. Bohl. The pack horse is not doing any work on the pack, because the pack is moving in a horizontal line at constant speed, and therefore there is no kinetic or gravitational energy being transferred into or out of it.

have

$$\begin{aligned} K &= \sum_{j=1}^n \frac{1}{2} m_j (\mathbf{v}_j + \mathbf{u}) \cdot (\mathbf{v}_j + \mathbf{u}) \\ &= \frac{1}{2} \left[\sum_{j=1}^n m_j \mathbf{v}_j \cdot \mathbf{v}_j + 2 \sum_{j=1}^n m_j \mathbf{v}_j \cdot \mathbf{u} + \sum_{j=1}^n m_j \mathbf{u} \cdot \mathbf{u} \right]. \end{aligned}$$

As in the proof on page 132, the first sum is simply the total kinetic energy in the original frame of reference, and the last sum is a constant, which has no effect on the validity of the conservation law. The middle sum can be rewritten as

$$\begin{aligned} 2 \sum_{j=1}^n m_j \mathbf{v}_j \cdot \mathbf{u} &= 2 \mathbf{u} \cdot \sum_{j=1}^n m_j \mathbf{v}_j \\ &= 2 \mathbf{u} \cdot \sum_{j=1}^n \mathbf{p}_j, \end{aligned}$$

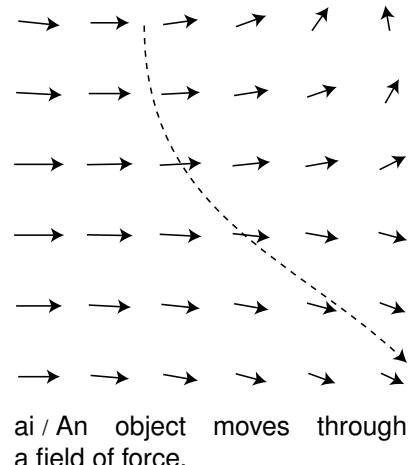
so the only way energy can be conserved for all values of \mathbf{u} is if the vector sum of the momenta is conserved as well.

3.4.6 Gradients and line integrals (optional)

This subsection introduces a little bit of vector calculus. It can be omitted without loss of continuity, but the techniques will be needed in our study of electricity and magnetism, and it may be helpful to be exposed to them in easy-to-visualize mechanical contexts before applying them to invisible electrical and magnetic phenomena.

In physics we often deal with fields of force, meaning situations where the force on an object depends on its position. For instance, figure ai could represent a map of the trade winds affecting a sailing ship, or a chart of the gravitational forces experienced by a space probe entering a double-star system. An object moving under the influence of this force will not necessarily be moving in the same direction as the force at every moment. The sailing ship can tack against the wind, due to the force from the water on the keel. The space probe, if it entered from the top of the diagram at high speed, would start to curve around to the right, but its inertia would carry it forward, and it wouldn't instantly swerve to match the direction of the gravitational force. For convenience, we've defined the gravitational field, \mathbf{g} , as the force *per unit mass*, but that trick only leads to a simplification because the gravitational force on an object is proportional to its mass. Since this subsection is meant to apply to any kind of force, we'll discuss everything in terms of the actual force vector, \mathbf{F} , in units of newtons.

If an object moves through the field of force along some curved path from point \mathbf{r}_1 to point \mathbf{r}_2 , the force will do a certain amount



of work on it. To calculate this work, we can break the path up into infinitesimally short segments, find the work done along each segment, and add them all up. For an object traveling along a nice straight x axis, we use the symbol dx to indicate the length of any infinitesimally short segment. In three dimensions, moving along a curve, each segment is a tiny vector $d\mathbf{r} = \hat{\mathbf{x}} dx + \hat{\mathbf{y}} dy + \hat{\mathbf{z}} dz$. The work theorem can be expressed as a dot product, so the work done along a segment is $\mathbf{F} \cdot d\mathbf{r}$. We want to integrate this, but we don't know how to integrate with respect to a variable that's a vector, so let's define a variable s that indicates the distance traveled so far along the curve, and integrate with respect to it instead. The expression $\mathbf{F} \cdot d\mathbf{r}$ can be rewritten as $|\mathbf{F}| |\mathbf{d}\mathbf{r}| \cos \theta$, where θ is the angle between \mathbf{F} and $d\mathbf{r}$. But $|\mathbf{d}\mathbf{r}|$ is simply ds , so the amount of work done becomes

$$\Delta E = \int_{\mathbf{r}_1}^{\mathbf{r}_2} |\mathbf{F}| \cos \theta \ ds.$$

Both \mathbf{F} and θ are functions of s . As a matter of notation, it's cumbersome to have to write the integral like this. Vector notation was designed to eliminate this kind of drudgery. We therefore define the line integral

$$\int_C \mathbf{F} \cdot d\mathbf{r}$$

as a way of notating this type of integral. The 'C' refers to the curve along which the object travels. If we don't know this curve then we typically can't evaluate the line integral just by knowing the initial and final positions \mathbf{r}_1 and \mathbf{r}_2 .

The basic idea of calculus is that integration undoes differentiation, and vice-versa. In one dimension, we could describe an interaction either in terms of a force or in terms of an interaction energy. We could integrate force with respect to position to find minus the energy, or we could find the force by taking minus the derivative of the energy. In the line integral, position is represented by a vector. What would it mean to take a derivative with respect to a vector? The correct way to generalize the derivative dU/dx to three dimensions is to replace it with the following vector,

$$\frac{dU}{dx} \hat{\mathbf{x}} + \frac{dU}{dy} \hat{\mathbf{y}} + \frac{dU}{dz} \hat{\mathbf{z}},$$

called the *gradient* of U , and written with an upside-down delta¹⁸ like this, ∇U . Each of these three derivatives is really what's known as a partial derivative. What that means is that when you're differentiating U with respect to x , you're supposed to treat y and z and constants, and similarly when you do the other two derivatives. To emphasize that a derivative is a partial derivative, it's customary to

¹⁸The symbol ∇ is called a "nabla." Cool word!

write it using the symbol ∂ in place of the differential d's. Putting all this notation together, we have

$$\nabla U = \frac{\partial U}{\partial x} \hat{\mathbf{x}} + \frac{\partial U}{\partial y} \hat{\mathbf{y}} + \frac{\partial U}{\partial z} \hat{\mathbf{z}} \quad [\text{definition of the gradient}].$$

The gradient looks scary, but it has a very simple physical interpretation. It's a vector that points in the direction in which U is increasing most rapidly, and it tells you how rapidly U is increasing in that direction. For instance, sperm cells in plants and animals find the egg cells by traveling in the direction of the gradient of the concentration of certain hormones. When they reach the location of the strongest hormone concentration, they find their destiny. In terms of the gradient, the force corresponding to a given interaction energy is $\mathbf{F} = -\nabla U$.

Force exerted by a spring

example 78

In one dimension, Hooke's law is $U = (1/2)kx^2$. Suppose we tether one end of a spring to a post, but it's free to stretch and swing around in a plane. Let's say its equilibrium length is zero, and let's choose the origin of our coordinate system to be at the post. Rotational invariance requires that its energy only depend on the magnitude of the \mathbf{r} vector, not its direction, so in two dimensions we have $U = (1/2)k|\mathbf{r}|^2 = (1/2)k(x^2 + y^2)$. The force exerted by the spring is then

$$\begin{aligned}\mathbf{F} &= -\nabla U \\ &= -\frac{\partial U}{\partial x} \hat{\mathbf{x}} - \frac{\partial U}{\partial y} \hat{\mathbf{y}} \\ &= -kx \hat{\mathbf{x}} - ky \hat{\mathbf{y}}.\end{aligned}$$

The magnitude of this force vector is $k|\mathbf{r}|$, and its direction is toward the origin.

This chapter is summarized on page 1077. Notation and terminology are tabulated on pages 1070-1071.

Problems

The symbols \checkmark , \blacksquare , etc. are explained on page 243.

1 Derive a formula expressing the kinetic energy of an object in terms of its momentum and mass. \checkmark \blacksquare

2 Two people in a rowboat wish to move around without causing the boat to move. What should be true about their total momentum? Explain. \blacksquare

3 A bullet leaves the barrel of a gun with a kinetic energy of 90 J. The gun barrel is 50 cm long. The gun has a mass of 4 kg, the bullet 10 g.

(a) Find the bullet's final velocity. \checkmark

(b) Find the bullet's final momentum. \checkmark

(c) Find the momentum of the recoiling gun.

(d) Find the kinetic energy of the recoiling gun, and explain why the recoiling gun does not kill the shooter. \checkmark \blacksquare

4 The big difference between the equations for momentum and kinetic energy is that one is proportional to v and one to v^2 . Both, however, are proportional to m . Suppose someone tells you that there's a third quantity, funkosity, defined as $f = m^2v$, and that funkosity is conserved. How do you know your leg is being pulled?

\triangleright Solution, p. 1040 \blacksquare

5 A ball of mass $2m$ collides head-on with an initially stationary ball of mass m . No kinetic energy is transformed into heat or sound. In what direction is the mass- $2m$ ball moving after the collision, and how fast is it going compared to its original velocity?

\triangleright Answer, p. 1068 \blacksquare

6 A very massive object with velocity v collides head-on with an object at rest whose mass is very small. No kinetic energy is converted into other forms. Prove that the low-mass object recoils with velocity $2v$. [Hint: Use the center-of-mass frame of reference.] \blacksquare

7 A mass m moving at velocity v collides with a stationary target having the same mass m . Find the maximum amount of energy that can be released as heat and sound. \checkmark \blacksquare

8 A rocket ejects exhaust with an exhaust velocity u . The rate at which the exhaust mass is used (mass per unit time) is b . We assume that the rocket accelerates in a straight line starting from rest, and that no external forces act on it. Let the rocket's initial mass (fuel plus the body and payload) be m_i , and m_f be its final mass, after all the fuel is used up. (a) Find the rocket's final velocity, v , in terms of u , m_i , and m_f . Neglect the effects of special relativity. (b) A typical exhaust velocity for chemical rocket engines is 4000 m/s. Estimate the initial mass of a rocket that could accelerate a one-ton payload to 10% of the speed of light, and show that this

design won't work. (For the sake of the estimate, ignore the mass of the fuel tanks. The speed is fairly small compared to c , so it's not an unreasonable approximation to ignore relativity.) \checkmark ■

9 An object is observed to be moving at constant speed along a line. Can you conclude that no forces are acting on it? Explain. [Based on a problem by Serway and Faughn.] ■

10 At low speeds, every car's acceleration is limited by traction, not by the engine's power. Suppose that at low speeds, a certain car is normally capable of an acceleration of 3 m/s^2 . If it is towing a trailer with half as much mass as the car itself, what acceleration can it achieve? [Based on a problem from PSSC Physics.] ■

11 (a) Let T be the maximum tension that an elevator's cable can withstand without breaking, i.e., the maximum force it can exert. If the motor is programmed to give the car an acceleration a ($a > 0$ is upward), what is the maximum mass that the car can have, including passengers, if the cable is not to break? \checkmark
(b) Interpret the equation you derived in the special cases of $a = 0$ and of a downward acceleration of magnitude g .

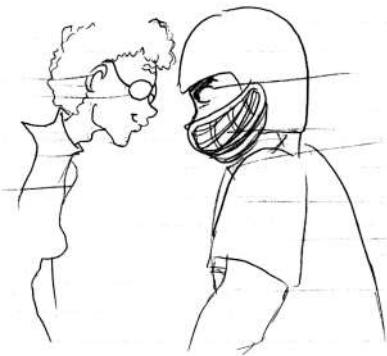


12 When the contents of a refrigerator cool down, the changed molecular speeds imply changes in both momentum and energy. Why, then, does a fridge transfer *power* through its radiator coils, but not *force*? \triangleright Solution, p. 1040 ■

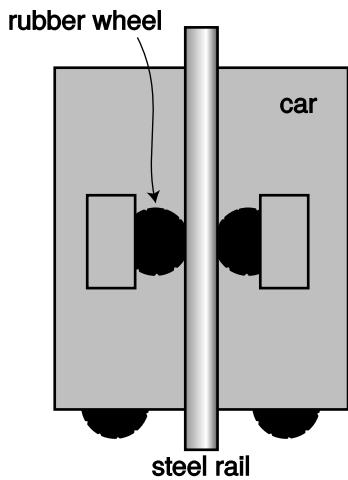
13 A helicopter of mass m is taking off vertically. The only forces acting on it are the earth's gravitational force and the force, F_{air} , of the air pushing up on the propeller blades.

- (a) If the helicopter lifts off at $t = 0$, what is its vertical speed at time t ?
- (b) Check that the units of your answer to part a make sense.
- (c) Discuss how your answer to part a depends on all three variables, and show that it makes sense. That is, for each variable, discuss what would happen to the result if you changed it while keeping the other two variables constant. Would a bigger value give a smaller result, or a bigger result? Once you've figured out this *mathematical* relationship, show that it makes sense *physically*.
- (d) Plug numbers into your equation from part a, using $m = 2300 \text{ kg}$, $F_{\text{air}} = 27000 \text{ N}$, and $t = 4.0 \text{ s}$. \checkmark ■

14 A blimp is initially at rest, hovering, when at $t = 0$ the pilot turns on the engine driving the propeller. The engine cannot instantly get the propeller going, but the propeller speeds up steadily. The steadily increasing force between the air and the propeller is given by the equation $F = kt$, where k is a constant. If the mass of the blimp is m , find its position as a function of time. (Assume that during the period of time you're dealing with, the blimp is not yet moving fast enough to cause a significant backward force due to



Problem 16.



Problem 19.

air resistance.) ✓ ■

15 A car is accelerating forward along a straight road. If the force of the road on the car's wheels, pushing it forward, is a constant 3.0 kN, and the car's mass is 1000 kg, then how long will the car take to go from 20 m/s to 50 m/s? ► Solution, p. 1040 ■

16 A little old lady and a pro football player collide head-on. Compare their forces on each other, and compare their accelerations. Explain. ■

17 The earth is attracted to an object with a force equal and opposite to the force of the earth on the object. If this is true, why is it that when you drop an object, the earth does not have an acceleration equal and opposite to that of the object? ■

18 When you stand still, there are two forces acting on you, the force of gravity (your weight) and the normal force of the floor pushing up on your feet. Are these forces equal and opposite? Does Newton's third law relate them to each other? Explain. ■

19 Today's tallest buildings are really not that much taller than the tallest buildings of the 1940's. One big problem with making an even taller skyscraper is that every elevator needs its own shaft running the whole height of the building. So many elevators are needed to serve the building's thousands of occupants that the elevator shafts start taking up too much of the space within the building. An alternative is to have elevators that can move both horizontally and vertically: with such a design, many elevator cars can share a few shafts, and they don't get in each other's way too much because they can detour around each other. In this design, it becomes impossible to hang the cars from cables, so they would instead have to ride on rails which they grab onto with wheels. Friction would keep them from slipping. The figure shows such a frictional elevator in its vertical travel mode. (The wheels on the bottom are for when it needs to switch to horizontal motion.)

(a) If the coefficient of static friction between rubber and steel is μ_s , and the maximum mass of the car plus its passengers is M , how much force must there be pressing each wheel against the rail in order to keep the car from slipping? (Assume the car is not accelerating.) ✓

(b) Show that your result has physically reasonable behavior with respect to μ_s . In other words, if there was less friction, would the wheels need to be pressed more firmly or less firmly? Does your equation behave that way? ■

20

A tugboat of mass m pulls a ship of mass M , accelerating it. Ignore fluid friction acting on their hulls, although there will of course need to be fluid friction acting on the tug's propellers.

- (a) If the force acting on the tug's propeller is F , what is the tension, T , in the cable connecting the two ships? \triangleright Hint, p. 1035 \checkmark
(b) Interpret your answer in the special cases of $M = 0$ and $M = \infty$.



21 Someone tells you she knows of a certain type of Central American earthworm whose skin, when rubbed on polished diamond, has $\mu_k > \mu_s$. Why is this not just empirically unlikely but logically suspect?



22 A uranium atom deep in the earth spits out an alpha particle. An alpha particle is a fragment of an atom. This alpha particle has initial speed v , and travels a distance d before stopping in the earth.

- (a) Find the force, F , from the dirt that stopped the particle, in terms of v, d , and its mass, m . Don't plug in any numbers yet. Assume that the force was constant. \checkmark

(b) Show that your answer has the right units.

(c) Discuss how your answer to part a depends on all three variables, and show that it makes sense. That is, for each variable, discuss what would happen to the result if you changed it while keeping the other two variables constant. Would a bigger value give a smaller result, or a bigger result? Once you've figured out this *mathematical* relationship, show that it makes sense *physically*.

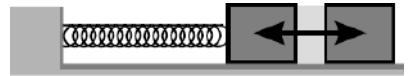
- (d) Evaluate your result for $m = 6.7 \times 10^{-27}$ kg, $v = 2.0 \times 10^4$ km/s, and $d = 0.71$ mm. \checkmark



23 You are given a large sealed box, and are not allowed to open it. Which of the following experiments measure its mass, and which measure its weight? [Hint: Which experiments would give different results on the moon?]

- (a) Put it on a frozen lake, throw a rock at it, and see how fast it scoots away after being hit.
(b) Drop it from a third-floor balcony, and measure how loud the sound is when it hits the ground.
(c) As shown in the figure, connect it with a spring to the wall, and watch it vibrate.

\triangleright Solution, p. 1040 \blacksquare



Problem 23, part c.

24 While escaping from the palace of the evil Martian emperor, Sally Spacehound jumps from a tower of height h down to the ground. Ordinarily the fall would be fatal, but she fires her blaster rifle straight down, producing an upward force of magnitude F_B . This force is insufficient to levitate her, but it does cancel out some of the force of gravity. During the time t that she is falling,

Sally is unfortunately exposed to fire from the emperor's minions, and can't dodge their shots. Let m be her mass, and g the strength of gravity on Mars.

- (a) Find the time t in terms of the other variables.
- (b) Check the units of your answer to part a.
- (c) For sufficiently large values of F_B , your answer to part a becomes nonsense — explain what's going on. ✓



25 When I cook rice, some of the dry grains always stick to the measuring cup. To get them out, I turn the measuring cup upside-down and hit the “roof” with my hand so that the grains come off of the “ceiling.” (a) Explain why static friction is irrelevant here. (b) Explain why gravity is negligible. (c) Explain why hitting the cup works, and why its success depends on hitting the cup hard enough.



26 A flexible rope of mass m and length L slides without friction over the edge of a table. Let x be the length of the rope that is hanging over the edge at a given moment in time.

- (a) Show that x satisfies the equation of motion $d^2x/dt^2 = gx/L$. [Hint: Use $F = dp/dt$, which allows you to handle the two parts of the rope separately even though mass is moving out of one part and into the other.]
- (b) Give a physical explanation for the fact that a larger value of x on the right-hand side of the equation leads to a greater value of the acceleration on the left side.
- (c) When we take the second derivative of the function $x(t)$ we are supposed to get essentially the same function back again, except for a constant out in front. The function e^x has the property that it is unchanged by differentiation, so it is reasonable to look for solutions to this problem that are of the form $x = be^{ct}$, where b and c are constants. Show that this does indeed provide a solution for two specific values of c (and for any value of b).
- (d) Show that the sum of any two solutions to the equation of motion is also a solution.
- (e) Find the solution for the case where the rope starts at rest at $t = 0$ with some nonzero value of x . ■

In problems 27-31, analyze the forces using a table in the format shown in section 3.2.6. Analyze the forces in which the italicized object participates.

27 Some people put a spare car key in a little magnetic *box* that they stick under the chassis of their car. Let's say that the box is stuck directly underneath a horizontal surface, and the car is parked. (See instructions above.) ■

28 Analyze two examples of *objects* at rest relative to the earth that are being kept from falling by forces other than the normal force. Do not use objects in outer space, and do not duplicate

problem 27 or 31. (See instructions above.)

29 A person is rowing a boat, with her feet braced. She is doing the part of the stroke that propels the boat, with the ends of the oars in the water (not the part where the oars are out of the water). (See instructions above.)

30 A farmer is in a stall with a cow when the cow decides to press him against the wall, pinning him with his feet off the ground. Analyze the forces in which the farmer participates. (See instructions above.)

31 A propeller *plane* is cruising east at constant speed and altitude. (See instructions above.)

32 The figure shows a stack of two blocks, sitting on top of a table that is bolted to the floor. All three objects are made from identical wood, with their surfaces finished identically using the same sandpaper. We tap the middle block, giving it an initial velocity v to the right. The tap is executed so rapidly that almost no initial velocity is imparted to the top block.

(a) Find the time that will elapse until the slipping between the top and middle blocks stops. Express your answer in terms of v , m , M , g , and the relevant coefficient of friction. ✓

(b) Show that your answer makes sense in terms of units.

(c) Check that your result has the correct behavior when you make m bigger or smaller. Explain. This means that you should discuss the mathematical behavior of the result, and then explain how this corresponds to what would really happen physically.

(d) Similarly, discuss what happens when you make M bigger or smaller.

(e) Similarly, discuss what happens when you make g bigger or smaller.

33 Ginny has a plan. She is going to ride her sled while her dog Foo pulls her, and she holds on to his leash. However, Ginny hasn't taken physics, so there may be a problem: she may slide right off the sled when Foo starts pulling.

(a) Analyze all the forces in which Ginny participates, making a table as in subsection 3.2.6.

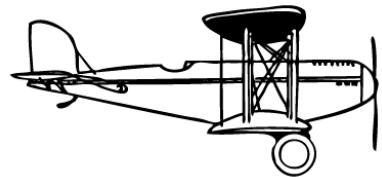
(b) Analyze all the forces in which the sled participates.

(c) The sled has mass m , and Ginny has mass M . The coefficient of static friction between the sled and the snow is μ_1 , and μ_2 is the corresponding quantity for static friction between the sled and her snow pants. Ginny must have a certain minimum mass so that she will not slip off the sled. Find this in terms of the other three variables. ✓

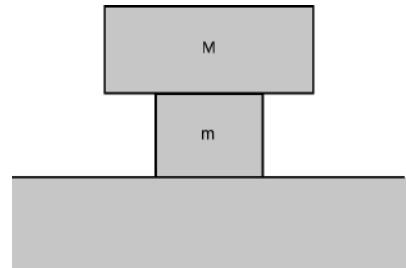
(d) Interpreting your equation from part c, under what conditions will there be no physically realistic solution for M ? Discuss what this means physically.



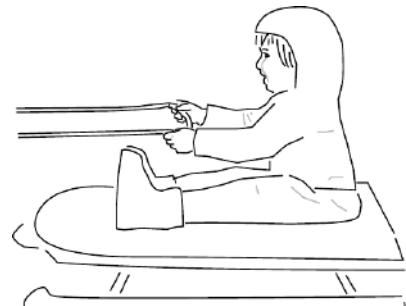
Problem 29.



Problem 31.



Problem 32



Problem 33.

34 In each case, identify the force that causes the acceleration, and give its Newton's-third-law partner. Describe the effect of the partner force. (a) A swimmer speeds up. (b) A golfer hits the ball off of the tee. (c) An archer fires an arrow. (d) A locomotive slows down.

► Solution, p. 1040 ■

35 A cop investigating the scene of an accident measures the length L of a car's skid marks in order to find out its speed v at the beginning of the skid. Express v in terms of L and any other relevant variables.

✓ ■

36 An ice skater builds up some speed, and then coasts across the ice passively in a straight line. (a) Analyze the forces, using a table in the format shown in subsection 3.2.6.

(b) If his initial speed is v , and the coefficient of kinetic friction is μ_k , find the maximum theoretical distance he can glide before coming to a stop. Ignore air resistance.

✓

(c) Show that your answer to part b has the right units.

(d) Show that your answer to part b depends on the variables in a way that makes sense physically.

(e) Evaluate your answer numerically for $\mu_k = 0.0046$, and a world-record speed of 14.58 m/s. (The coefficient of friction was measured by De Koning et al., using special skates worn by real speed skaters.)

✓

(f) Comment on whether your answer in part e seems realistic. If it doesn't, suggest possible reasons why.

■

37 (a) Using the solution of problem 37 on page 126, predict how the spring constant of a fiber will depend on its length and cross-sectional area.

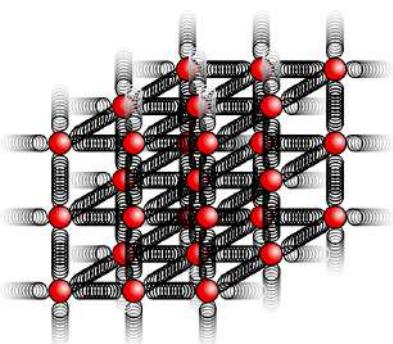
(b) The constant of proportionality is called the Young's modulus, E , and typical values of the Young's modulus are about 10^{10} to 10^{11} . What units would the Young's modulus have in the SI system?

► Solution, p. 1040 ■

38 This problem depends on the results of problems problem 37 on page 126 and problem 37 from this chapter. When atoms form chemical bonds, it makes sense to talk about the spring constant of the bond as a measure of how "stiff" it is. Of course, there aren't really little springs — this is just a mechanical model. The purpose of this problem is to estimate the spring constant, k , for a single bond in a typical piece of solid matter. Suppose we have a fiber, like a hair or a piece of fishing line, and imagine for simplicity that it is made of atoms of a single element stacked in a cubical manner, as shown in the figure, with a center-to-center spacing b . A typical value for b would be about 10^{-10} m.

(a) Find an equation for k in terms of b , and in terms of the Young's modulus, E , defined in problem 37 and its solution.

(b) Estimate k using the numerical data given in problem 37.



Problem 38

(c) Suppose you could grab one of the atoms in a diatomic molecule like H₂ or O₂, and let the other atom hang vertically below it. Does the bond stretch by any appreciable fraction due to gravity? ■

39 This problem has been deleted. ■

40 Many fish have an organ known as a swim bladder, an air-filled cavity whose main purpose is to control the fish's buoyancy and allow it to keep from rising or sinking without having to use its muscles. In some fish, however, the swim bladder (or a small extension of it) is linked to the ear and serves the additional purpose of amplifying sound waves. For a typical fish having such an anatomy, the bladder has a resonant frequency of 300 Hz, the bladder's Q is 3, and the maximum amplification is about a factor of 100 in energy. Over what range of frequencies would the amplification be at least a factor of 50?

✓ ■

41 An oscillator with sufficiently strong damping has its maximum response at $\omega = 0$. Using the result derived on page 1027 , find the value of Q at which this behavior sets in.

▷ Hint, p. 1035 ▷ Answer, p. 1068 ■

42 An oscillator has $Q=6.00$, and, for convenience, let's assume $F_m = 1.00$, $\omega_0 = 1.00$, and $m = 1.00$. The usual approximations would give

$$\begin{aligned}\omega_{res} &= \omega_0, \\ A_{res} &= 6.00, \quad \text{and} \\ \Delta\omega &= 1/6.00.\end{aligned}$$

Determine these three quantities numerically using the result derived on page 1027 , and compare with the approximations. ■

43 The apparatus in figure d on page 61 had a natural period of oscillation of 5 hours and 20 minutes. The authors estimated, based on calculations of internal friction in the tungsten wire, that its Q was on the order of 10^6 , but they were unable to measure it empirically because it would have taken years for the amplitude to die down by any measurable amount. Although each aluminum or platinum mass was really moving along an arc of a circle, any actual oscillations caused by a violation of the equivalence of gravitational and inertial mass would have been measured in millions of a degree, so it's a good approximation to say that each mass's motion was along a (very short!) straight line segment. We can also treat each

mass as if it was oscillating separately from the others. If the principle of equivalence had been violated at the 10^{-12} level, the limit of their experiment's sensitivity, the sun's gravitational force on one of the 0.4-gram masses would have been about 3×10^{-19} N, oscillating with a period of 24 hours due to the rotation of the earth. (We ignore the inertia of the arms, whose total mass was only about 25% of the total mass of the rotating assembly.)

(a) Find the amplitude of the resulting oscillations, and determine the angle to which they would have corresponded, given that the radius of the balance arms was 10 cm. \triangleright Answer, p. 1068

(b) Show that even if their estimate of Q was wildly wrong, it wouldn't have affected this result. \blacksquare

44 A firework shoots up into the air, and just before it explodes it has a certain momentum and kinetic energy. What can you say about the momenta and kinetic energies of the pieces immediately after the explosion? [Based on a problem from PSSC Physics.]

\triangleright Solution, p. 1040 \blacksquare

45 The figure shows a view from above of a collision about to happen between two air hockey pucks sliding without friction. They have the same speed, v_i , before the collision, but the big puck is 2.3 times more massive than the small one. Their sides have sticky stuff on them, so when they collide, they will stick together. At what angle will they emerge from the collision? In addition to giving a numerical answer, please indicate by drawing on the figure how your angle is defined. \triangleright Solution, p. 1041 \blacksquare



Problem 45

46 A learjet traveling due east at 300 mi/hr collides with a jumbo jet which was heading southwest at 150 mi/hr. The jumbo jet's mass is five times greater than that of the learjet. When they collide, the learjet sticks into the fuselage of the jumbo jet, and they fall to earth together. Their engines stop functioning immediately after the collision. On a map, what will be the direction from the location of the collision to the place where the wreckage hits the ground? (Give an angle.) \checkmark \blacksquare

47 (a) A ball is thrown straight up with velocity v . Find an equation for the height to which it rises. \checkmark

(b) Generalize your equation for a ball thrown at an angle θ above horizontal, in which case its initial velocity components are $v_x = v \cos \theta$ and $v_y = v \sin \theta$. \checkmark \blacksquare

48 At the 2010 Salinas Lettuce Festival Parade, the Lettuce Queen drops her bouquet while riding on a float moving toward the right. Sketch the shape of its trajectory in her frame of reference, and compare with the shape seen by one of her admirers standing on the sidewalk. \blacksquare

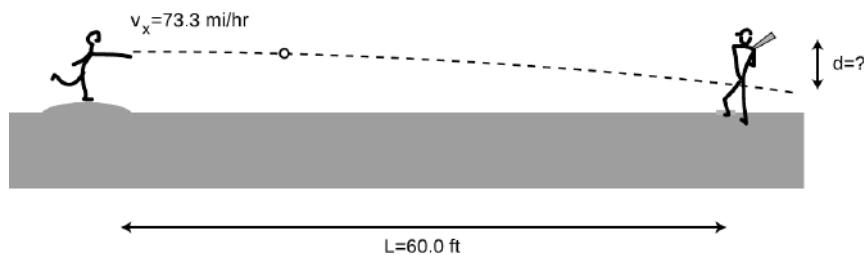
49 Two daredevils, Wendy and Bill, go over Niagara Falls.

Wendy sits in an inner tube, and lets the 30 km/hr velocity of the river throw her out horizontally over the falls. Bill paddles a kayak, adding an extra 10 km/hr to his velocity. They go over the edge of the falls at the same moment, side by side. Ignore air friction. Explain your reasoning.

- (a) Who hits the bottom first?
- (b) What is the horizontal component of Wendy's velocity on impact?
- (c) What is the horizontal component of Bill's velocity on impact?
- (d) Who is going faster on impact? ■

50 A baseball pitcher throws a pitch clocked at $v_x = 73.3$ miles/hour. He throws horizontally. By what amount, d , does the ball drop by the time it reaches home plate, $L = 60.0$ feet away?

- (a) First find a symbolic answer in terms of L , v_x , and g . ✓
- (b) Plug in and find a numerical answer. Express your answer in units of ft. (Note: 1 foot=12 inches, 1 mile=5280 feet, and 1 inch=2.54 cm) ✓ ■



Problem 50.

51 A batter hits a baseball at speed v , at an angle θ above horizontal.

- (a) Find an equation for the range (horizontal distance to where the ball falls), R , in terms of the relevant variables. Neglect air friction and the height of the ball above the ground when it is hit.
▷ Answer, p. 1068
- (b) Interpret your equation in the cases of $\theta=0$ and $\theta = 90^\circ$.
- (c) Find the angle that gives the maximum range.
▷ Answer, p. 1068 ■

52 In this problem you'll extend the analysis in problem 51

to include air friction by writing a computer program. For a game played at sea level, the force due to air friction is approximately $(7 \times 10^{-4} \text{ N}\cdot\text{s}^2/\text{m}^2)v^2$, in the direction opposite to the motion of the ball. The mass of a baseball is 0.146 kg.

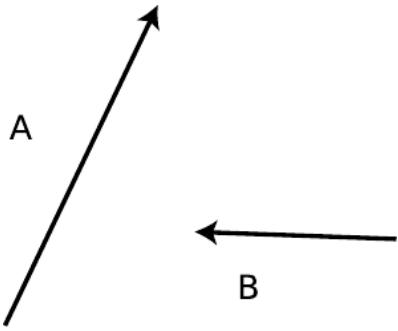
- (a) For a ball hit at a speed of 45.0 m/s from a height of 1.0 m, find the optimal angle and the resulting range. ▷ Answer, p. 1068
- (b) How much farther would the ball fly at the Colorado Rockies' stadium, where the thinner air gives 18 percent less air friction?

▷ Answer, p. 1068 ■

53

If you walk 35 km at an angle 25° counterclockwise from east, and then 22 km at 230° counterclockwise from east, find the distance and direction from your starting point to your destination.

✓ ■



Problem 54.

54 The figure shows vectors **A** and **B**. Graphically calculate the following, as in figure q on p. 204.

$$\mathbf{A} + \mathbf{B}, \mathbf{A} - \mathbf{B}, \mathbf{B} - \mathbf{A}, -2\mathbf{B}, \mathbf{A} - 2\mathbf{B}$$

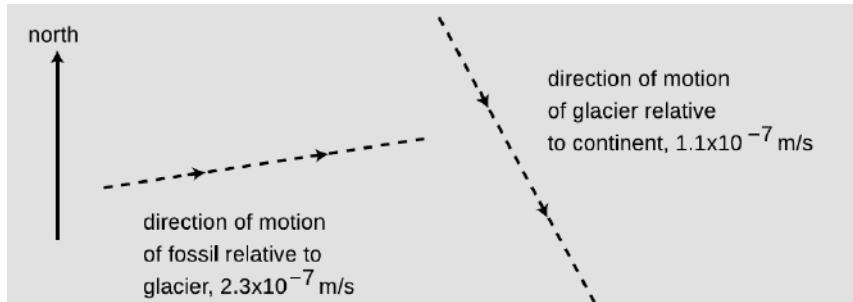
No numbers are involved. ■

55 Phnom Penh is 470 km east and 250 km south of Bangkok. Hanoi is 60 km east and 1030 km north of Phnom Penh.

(a) Choose a coordinate system, and translate these data into Δx and Δy values with the proper plus and minus signs.

(b) Find the components of the $\Delta\mathbf{r}$ vector pointing from Bangkok to Hanoi. ✓ ■

56 Is it possible for a helicopter to have an acceleration due east and a velocity due west? If so, what would be going on? If not, why not? ■

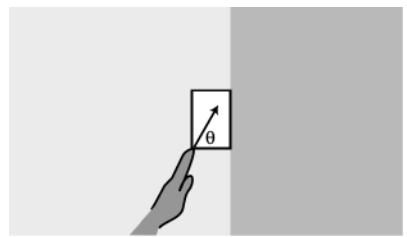


Problem 57.

57 As shown in the diagram, a dinosaur fossil is slowly moving down the slope of a glacier under the influence of wind, rain and gravity. At the same time, the glacier is moving relative to the continent underneath. The dashed lines represent the directions but not the magnitudes of the velocities. Pick a scale, and use graphical addition of vectors to find the magnitude and the direction of the fossil's velocity relative to the continent. You will need a ruler and protractor. ✓ ■

58 A bird is initially flying horizontally east at 21.1 m/s, but one second later it has changed direction so that it is flying horizontally and 7° north of east, at the same speed. What are the magnitude and direction of its acceleration vector during that one second time interval? (Assume its acceleration was roughly constant.) ✓ ■

59 Your hand presses a block of mass m against a wall with a force \mathbf{F}_H acting at an angle θ , as shown in the figure. Find the minimum and maximum possible values of $|\mathbf{F}_H|$ that can keep the block stationary, in terms of m , g , θ , and μ_s , the coefficient of static friction between the block and the wall. Check both your answers in the case of $\theta = 90^\circ$, and interpret the case where the maximum force is infinite. ✓ ■



Problem 59

60 A skier of mass m is coasting down a slope inclined at an angle θ compared to horizontal. Assume for simplicity that the treatment of kinetic friction given in chapter 5 is appropriate here, although a soft and wet surface actually behaves a little differently. The coefficient of kinetic friction acting between the skis and the snow is μ_k , and in addition the skier experiences an air friction force of magnitude bv^2 , where b is a constant.

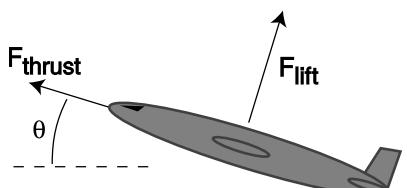
- (a) Find the maximum speed that the skier will attain, in terms of the variables m , g , θ , μ_k , and b . ✓
- (b) For angles below a certain minimum angle θ_{min} , the equation gives a result that is not mathematically meaningful. Find an equation for θ_{min} , and give a physical explanation of what is happening for $\theta < \theta_{min}$. ✓ ■

61 A gun is aimed horizontally to the west. The gun is fired, and the bullet leaves the muzzle at $t = 0$. The bullet's position vector as a function of time is $\mathbf{r} = b\hat{\mathbf{x}} + ct\hat{\mathbf{y}} + dt^2\hat{\mathbf{z}}$, where b , c , and d are positive constants.

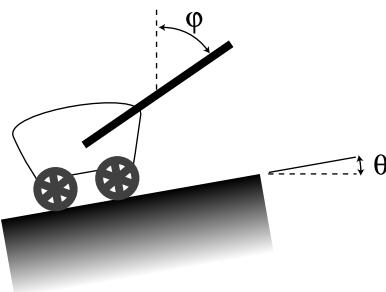
- (a) What units would b , c , and d need to have for the equation to make sense?
- (b) Find the bullet's velocity and acceleration as functions of time.
- (c) Give physical interpretations of b , c , d , $\hat{\mathbf{x}}$, $\hat{\mathbf{y}}$, and $\hat{\mathbf{z}}$. ■

62 Annie Oakley, riding north on horseback at 30 mi/hr, shoots her rifle, aiming horizontally and to the northeast. The muzzle speed of the rifle is 140 mi/hr. When the bullet hits a defenseless fuzzy animal, what is its speed of impact? Neglect air resistance, and ignore the vertical motion of the bullet. ▷ Solution, p. 1041 ■

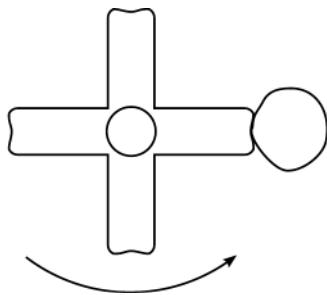
63 A cargo plane has taken off from a tiny airstrip in the Andes, and is climbing at constant speed, at an angle of $\theta = 17^\circ$ with respect to horizontal. Its engines supply a thrust of $F_{thrust} = 200$ kN, and the lift from its wings is $F_{lift} = 654$ kN. Assume that air resistance (drag) is negligible, so the only forces acting are thrust, lift, and weight. What is its mass, in kg? ▷ Solution, p. 1041 ■



Problem 63



Problem 64



Problem 66.

- 64** A wagon is being pulled at constant speed up a slope θ by a rope that makes an angle ϕ with the vertical. (a) Assuming negligible friction, show that the tension in the rope is given by the equation

$$T = \frac{\sin \theta}{\sin(\theta + \phi)} mg,$$

- (b) Interpret this equation in the special cases of $\phi = 0$ and $\phi = 180^\circ - \theta$. ▷ Solution, p. 1042 ■

- 65** The angle of repose is the maximum slope on which an object will not slide. On airless, geologically inert bodies like the moon or an asteroid, the only thing that determines whether dust or rubble will stay on a slope is whether the slope is less steep than the angle of repose.

- (a) Find an equation for the angle of repose, deciding for yourself what are the relevant variables.
 (b) On an asteroid, where g can be thousands of times lower than on Earth, would rubble be able to lie at a steeper angle of repose?
▷ Solution, p. 1042 ■

- 66** When you're done using an electric mixer, you can get most of the batter off of the beaters by lifting them out of the batter with the motor running at a high enough speed. Let's imagine, to make things easier to visualize, that we instead have a piece of tape stuck to one of the beaters.

- (a) Explain why static friction has no effect on whether or not the tape flies off.
 (b) Analyze the forces in which the tape participates, using a table in the format shown in subsection 3.2.6.
 (c) Suppose you find that the tape doesn't fly off when the motor is on a low speed, but at a greater speed, the tape won't stay on. Why would the greater speed change things? [Hint: If you don't invoke any law of physics, you haven't explained it.] ■

- 67** Show that the expression $|\mathbf{v}|^2/r$ has the units of acceleration. ■

- 68** A plane is flown in a loop-the-loop of radius 1.00 km. The plane starts out flying upside-down, straight and level, then begins curving up along the circular loop, and is right-side up when it reaches the top. (The plane may slow down somewhat on the way up.) How fast must the plane be going at the top if the pilot is to experience no force from the seat or the seatbelt while at the top of the loop? ✓ ■

- 69** Find the angle between the following two vectors:

$$\hat{\mathbf{x}} + 2\hat{\mathbf{y}} + 3\hat{\mathbf{z}}$$

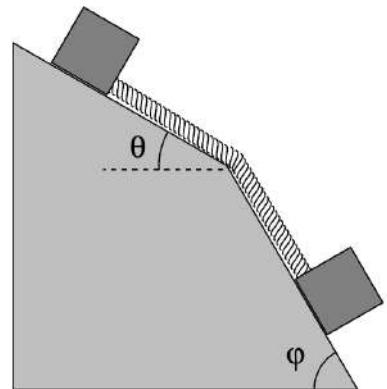
$$4\hat{\mathbf{x}} + 5\hat{\mathbf{y}} + 6\hat{\mathbf{z}}$$

▷ Hint, p. 1035 ✓ ■

70 The two blocks shown in the figure have equal mass, m , and the surface is frictionless. (a) What is the tension in the massless rope? \triangleright Hint, p. 1035 ✓

(b) Show that the units of your answer make sense.

(c) Check the physical behavior of your answer in the special cases of $\phi \leq \theta$ and $\theta = 0, \phi = 90^\circ$. ■



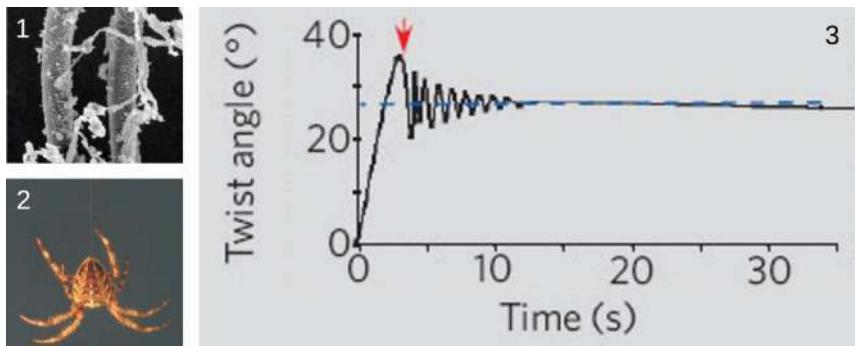
Problem 70.

71 (a) We observe that the amplitude of a certain free oscillation decreases from A_0 to A_0/Z after n oscillations. Find its Q . ✓

(b) The figure is from *Shape memory in Spider draglines*, Emile, Le Floch, and Vollrath, *Nature* 440:621 (2006). Panel 1 shows an electron microscope's image of a thread of spider silk. In 2, a spider is hanging from such a thread. From an evolutionary point of view, it's probably a bad thing for the spider if it twists back and forth while hanging like this. (We're referring to a back-and-forth rotation about the axis of the thread, not a swinging motion like a pendulum.) The authors speculate that such a vibration could make the spider easier for predators to see, and it also seems to me that it would be a bad thing just because the spider wouldn't be able to control its orientation and do what it was trying to do. Panel 3 shows a graph of such an oscillation, which the authors measured using a video camera and a computer, with a 0.1 g mass hung from it in place of a spider. Compared to human-made fibers such as kevlar or copper wire, the spider thread has an unusual set of properties:

1. It has a low Q , so the vibrations damp out quickly.
2. It doesn't become brittle with repeated twisting as a copper wire would.
3. When twisted, it tends to settle in to a new equilibrium angle, rather than insisting on returning to its original angle. You can see this in panel 2, because although the experimenters initially twisted the wire by 35 degrees, the thread only performed oscillations with an amplitude much smaller than ± 35 degrees, settling down to a new equilibrium at 27 degrees.
4. Over much longer time scales (hours), the thread eventually resets itself to its original equilibrium angle (shown as zero degrees on the graph). (The graph reproduced here only shows the motion over a much shorter time scale.) Some human-made materials have this "memory" property as well, but they typically need to be heated in order to make them go back to their original shapes.

Focusing on property number 1, estimate the Q of spider silk from the graph. ✓ ■



Problem 71.

72 A cross-country skier is gliding on a level trail, with negligible friction. Then, when he is at position $x = 0$, the tip of his skis enters a patch of dirt. As he rides onto the dirt, more and more of his weight is being supported by the dirt. The skis have length ℓ , so if he reached $x = \ell$ without stopping, his weight would be completely on the dirt. This problem deals with the motion for $x < \ell$.
 (a) Find the acceleration in terms of x , as well as any other relevant constants.

(b) This is a second-order differential equation. You should be able to find the solution simply by thinking about some commonly occurring functions that you know about, and finding two that have the right properties. If these functions are $x = f(t)$ and $x = g(t)$, then the most general solution to the equations of motion will be of the form $x = af + bg$, where a and b are constants to be determined from the initial conditions.

(c) Suppose that the initial velocity v_0 at $x = 0$ is such that he stops at $x < \ell$. Find the time until he stops, and show that, counterintuitively, this time is independent of v_0 . Explain physically why this is true. ✓

73 A microwave oven works by twisting molecules one way and then the other, counterclockwise and then clockwise about their own centers, millions of times a second. If you put an ice cube or a stick of butter in a microwave, you'll observe that the solid doesn't heat very quickly, although eventually melting begins in one small spot. Once this spot forms, it grows rapidly, while the rest of the solid remains solid; it appears that a microwave oven heats a liquid much more rapidly than a solid. Explain why this should happen, based on the atomic-level description of heat, solids, and liquids. (See, e.g., figure b on page 110.)

Don't repeat the following common mistakes:

In a solid, the atoms are packed more tightly and have less space between them. Not true. Ice floats because it's less dense than water.

In a liquid, the atoms are moving much faster. No, the difference in average speed between ice at -1°C and water at 1°C is only 0.4%. ■

74 Problem 2-16 on page 122 was intended to be solved using conservation of energy. Solve the same problem using Newton's laws. ■

75 A bead slides down along a piece of wire that is in the shape of a helix. The helix lies on the surface of a vertical cylinder of radius r , and the vertical distance between turns is d .

(a) Ordinarily when an object slides downhill under the influence of kinetic friction, the velocity-independence of kinetic friction implies that the acceleration is constant, and therefore there is no limit to the object's velocity. Explain the physical reason why this argument fails here, so that the bead will in fact have some limiting velocity.

(b) Find the limiting velocity.

(c) Show that your result has the correct behavior in the limit of $r \rightarrow \infty$. [Problem by B. Korsunsky.] ✓ ■

76 A person on a bicycle is to coast down a ramp of height h and then pass through a circular loop of radius r . What is the smallest value of h for which the cyclist will complete the loop without falling? (Ignore the kinetic energy of the spinning wheels.) ✓ ■

77 A car accelerates from rest. At low speeds, its acceleration is limited by static friction, so that if we press too hard on the gas, we will "burn rubber" (or, for many newer cars, a computerized traction-control system will override the gas pedal). At higher speeds, the limit on acceleration comes from the power of the engine, which puts a limit on how fast kinetic energy can be developed.

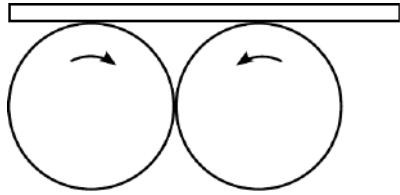
(a) Show that if a force F is applied to an object moving at speed v , the power required is given by $P = vF$.

(b) Find the speed v at which we cross over from the first regime described above to the second. At speeds higher than this, the engine does not have enough power to burn rubber. Express your result in terms of the car's power P , its mass m , the coefficient of static friction μ_s , and g . ✓

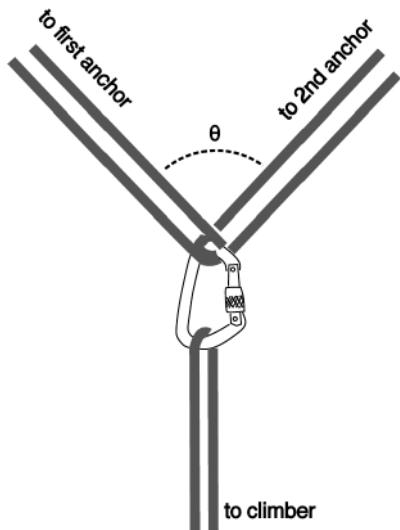
(c) Show that your answer to part b has units that make sense.

(d) Show that the dependence of your answer on each of the four variables makes sense physically.

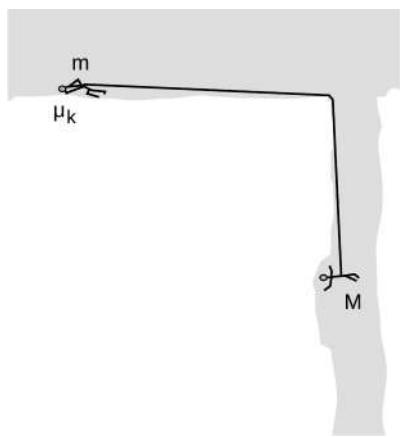
(e) The 2010 Maserati Gran Turismo Convertible has a maximum power of $3.23 \times 10^5 \text{ W}$ (433 horsepower) and a mass (including a 50-kg driver) of $2.03 \times 10^3 \text{ kg}$. (This power is the maximum the engine can supply at its optimum frequency of 7600 r.p.m. Presumably the automatic transmission is designed so a gear is available in which the engine will be running at very nearly this frequency when the car is moving at v .) Rubber on asphalt has $\mu_s \approx 0.9$. Find v for this car. Answer: 18 m/s, or about 40 miles per hour.



Problem 78.



Problem 79.



Problem 80.

(f) Our analysis has neglected air friction, which can probably be approximated as a force proportional to v^2 . The existence of this force is the reason that the car has a maximum speed, which is 176 miles per hour. To get a feeling for how good an approximation it is to ignore air friction, find what fraction of the engine's maximum power is being used to overcome air resistance when the car is moving at the speed v found in part e. Answer: 1% ■

78 Two wheels of radius r rotate in the same vertical plane with angular velocities $+\Omega$ and $-\Omega$ (rates of rotation in radians per second) about axes that are parallel and at the same height. The wheels touch one another at a point on their circumferences, so that their rotations mesh like gears in a gear train. A board is laid on top of the wheels, so that two friction forces act upon it, one from each wheel. Characterize the three qualitatively different types of motion that the board can exhibit, depending on the initial conditions. ■

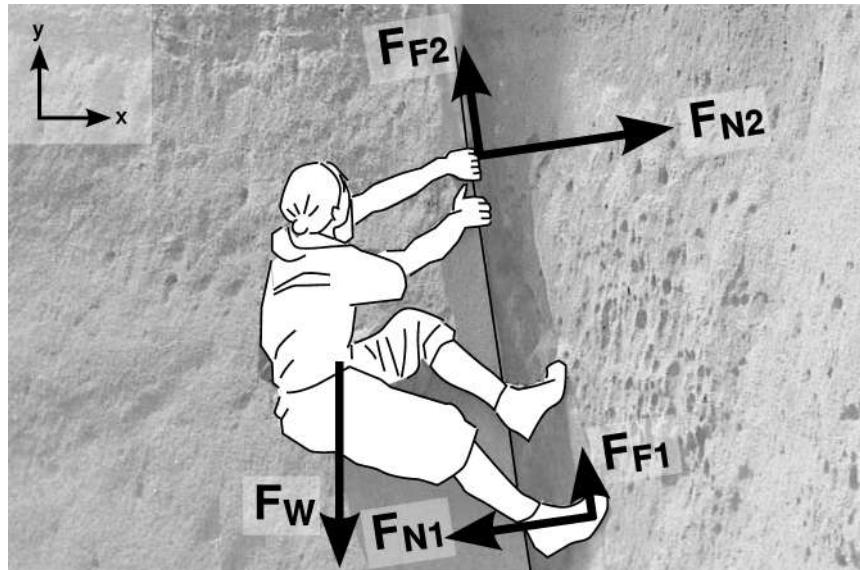
79 For safety, mountain climbers often wear a climbing harness and tie in to other climbers on a rope team or to anchors such as pitons or snow anchors. When using anchors, the climber usually wants to tie in to more than one, both for extra strength and for redundancy in case one fails. The figure shows such an arrangement, with the climber hanging from a pair of anchors forming a symmetric "Y" at an angle θ . The metal piece at the center is called a carabiner. The usual advice is to make $\theta < 90^\circ$; for large values of θ , the stress placed on the anchors can be many times greater than the actual load L , so that two anchors are actually *less* safe than one.

- (a) Find the force S at each anchor in terms of L and θ . ✓
- (b) Verify that your answer makes sense in the case of $\theta = 0$.
- (c) Interpret your answer in the case of $\theta = 180^\circ$.
- (d) What is the smallest value of θ for which S equals or exceeds L , so that for larger angles a failure of at least one anchor is *more* likely than it would have been with a single anchor? ✓

80 Mountain climbers with masses m and M are roped together while crossing a horizontal glacier when a vertical crevasse opens up under the climber with mass M . The climber with mass m drops down on the snow and tries to stop by digging into the snow with the pick of an ice ax. Alas, this story does not have a happy ending, because this doesn't provide enough friction to stop. Both m and M continue accelerating, with M dropping down into the crevasse and m being dragged across the snow, slowed only by the kinetic friction with coefficient μ_k acting between the ax and the snow. There is no significant friction between the rope and the lip of the crevasse.

- (a) Find the acceleration a . ✓
- (b) Check the units of your result.
- (c) Check the dependence of your equation on the variables. That means that for each variable, you should determine what its effect

on a should be physically, and then what your answer from part a says its effect would be mathematically. ■

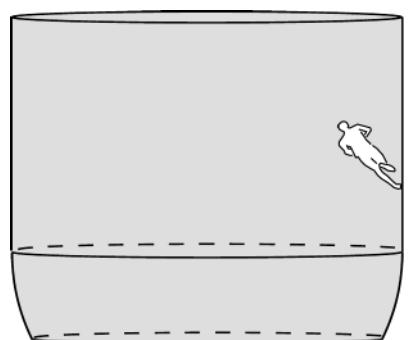


Problem 81.

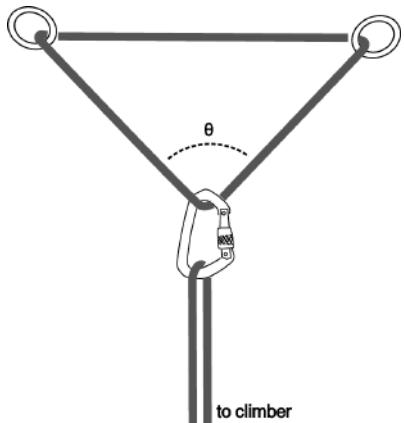
81 Complete example 71 on p. 209 by expressing the remaining nine x and y components of the forces in terms of the five magnitudes and the small, positive angle $\theta \approx 9^\circ$ by which the crack overhangs. ✓ ■

82 In a well known stunt from circuses and carnivals, a motorcyclist rides around inside a big bowl, gradually speeding up and rising higher. Eventually the cyclist can get up to where the walls of the bowl are vertical. Let's estimate the conditions under which a running human could do the same thing.

- (a) If the runner can run at speed v , and her shoes have a coefficient of static friction μ_s , what is the maximum radius of the circle? ✓
- (b) Show that the units of your answer make sense.
- (c) Check that its dependence on the variables makes sense.
- (d) Evaluate your result numerically for $v = 10$ m/s (the speed of an olympic sprinter) and $\mu_s = 5$. (This is roughly the highest coefficient of static friction ever achieved for surfaces that are not sticky. The surface has an array of microscopic fibers like a hair brush, and is inspired by the hairs on the feet of a gecko. These assumptions are not necessarily realistic, since the person would have to run at an angle, which would be physically awkward.) ✓ ■



Problem 82.



Problem 83.

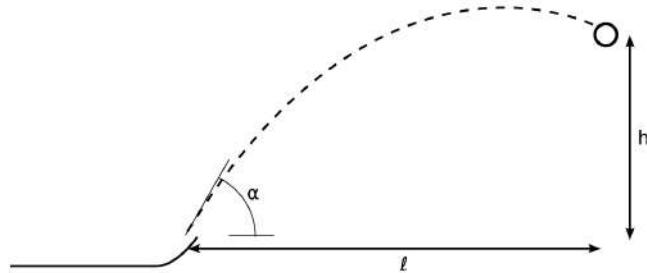
- 83** Problem 79 discussed a possible correct way of setting up a redundant anchor for mountaineering. The figure for this problem shows an incorrect way of doing it, by arranging the rope in a triangle (which we'll take to be isosceles). One of the bad things about the triangular arrangement is that it requires more force from the anchors, making them more likely to fail. (a) Using the same notation as in problem 79, find S in terms of L and θ . ✓
 (b) Verify that your answer makes sense in the case of $\theta = 0$, and compare with the correct setup. ■

- 84** At a picnic, someone hands you a can of beer. The ground is uneven, and you don't want to spill your drink. You reason that it will be more stable if you drink some of it first in order to lower its center of mass. How much should you drink in order to make the center of mass as low as possible? [Based on a problem by Walter van B. Roberts and Martin Gardner.] ■

- 85** “Big wall” climbing is a specialized type of rock climbing that involves going up tall cliffs such as the ones in Yosemite, usually with the climbers spending at least one night sleeping on a natural ledge or an artificial “portaledge.” In this style of climbing, each pitch of the climb involves strenuously hauling up several heavy bags of gear — a fact that has caused these climbs to be referred to as “vertical ditch digging.” (a) If an 80 kg haul bag has to be pulled up the full length of a 60 m rope, how much work is done? (b) Since it can be difficult to lift 80 kg, a 2:1 pulley is often used. The hauler then lifts the equivalent of 40 kg, but has to pull in 120 m of rope. How much work is done in this case? ✓ ■

86 The figure shows an arcade game called skee ball that is similar to bowling. The player rolls the ball down a horizontal alley. The ball then rides up a curved lip and is launched at an initial speed u , at an angle α above horizontal. Suppose we want the ball to go into a hole that is at horizontal distance ℓ and height h , as shown in the figure.

- (a) Find the initial speed u that is required, in terms of the other variables and g . ✓
- (b) Check that your answer to part a has units that make sense.
- (c) Check that your answer to part a depends on g in a way that makes sense. This means that you should first determine on physical grounds whether increasing g should increase u , or decrease it. Then see whether your answer to part a has this mathematical behavior.
- (d) Do the same for the dependence on h .
- (e) Interpret your equation in the case where $\alpha = 90^\circ$.
- (f) Interpret your equation in the case where $\tan \alpha = h/\ell$.
- (g) Find u numerically if $h = 70$ cm, $\ell = 60$ cm, and $\alpha = 65^\circ$. ✓



Problem 86.

87 The figure shows the International Space Station (ISS). One of the purposes of the ISS is supposed to be to carry out experiments in microgravity. However, the following factor limits this application. The ISS orbits the earth once every 92.6 minutes. It is desirable to keep the same side of the station always oriented toward the earth, which means that the station has to rotate with the same period. In the photo, the direction of orbital motion is left or right on the page, so the rotation is about the axis shown as up and down on the page. The greatest distance of any pressurized compartment from the axis of rotation is 36.5 meters. Find the acceleration due to the rotation at this point, and the apparent weight of a 60 kg astronaut at that location. ✓



Problem 87.

88 Problems 88-90 all investigate the following idea. Cosmological surveys at the largest observable distance scales have detected structures like filaments. As an idealization of such a structure, consider a uniform mass distribution lying along the entire x axis, with mass density λ in units of kg/m. The purpose of this problem is to find the gravitational field created by this structure at a distance y .

(a) Determine as much as possible about the form of the solution, based on units.

(b) To evaluate the actual result, find the contribution dg_y to the y component of the field arising from the mass dm lying between x and $x + dx$, then integrate it. \triangleright Solution, p. 1043 ■

89 Let us slightly change the physical situation described in problem 88, letting the filament have a finite size, while retaining its symmetry under rotation about the x axis. The details don't actually matter very much for our purposes, but if we like, we can take the mass density to be constant within a cylinder of radius b centered on the x axis. Now consider the following two limits:

$$g_1 = \lim_{y \rightarrow 0} \lim_{b \rightarrow 0} g \quad \text{and}$$

$$g_2 = \lim_{b \rightarrow 0} \lim_{y \rightarrow 0} g.$$

Each of these is a limit inside another limit, the only difference being the order of the limits. Either of these could be used as a definition of the field at a point *on* an infinitely thin filament. Do they agree? ■

90 Suppose we have a mass filament like the one described in problems 88 and 89, but now rather than taking it to be straight, let it have the shape of an arbitrary smooth curve. Locally, "under a microscope," this curve will look like an arc of a circle, i.e., we can describe its shape solely in terms of a radius of curvature. As in problem 89, consider a point P lying *on* the filament itself, taking g to be defined as in definition g_1 . Investigate whether g is finite, and also whether it points in a specific direction. To clarify the mathematical idea, consider the following two limits:

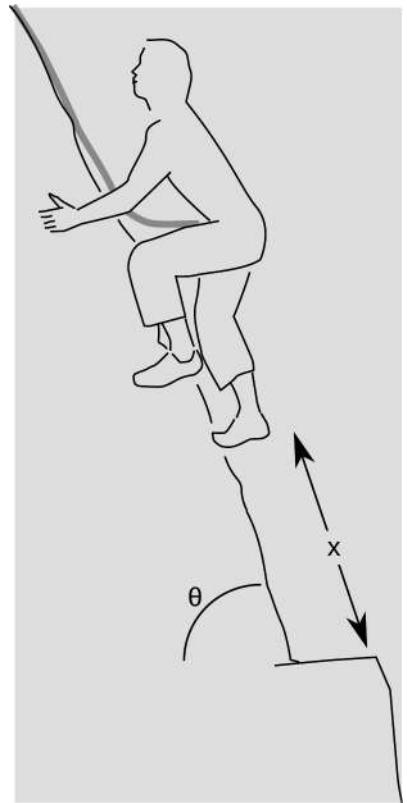
$$A = \lim_{x \rightarrow 0} \frac{1}{x} \quad \text{and}$$

$$B = \lim_{x \rightarrow 0} \frac{1}{x^2}.$$

We say that $A = \infty$, while $B = +\infty$, i.e., both diverge, but B diverges with a definite sign. For a straight filament, as in problem 88, with an infinite radius of curvature, symmetry guarantees that the field at P has no specific direction, in analogy with limit A. For a curved filament, a calculation is required in order to determine whether we get behavior A or B. Based on your result, what is the expected dynamical behavior of such a filament? ■

91 The rock climber in the figure has mass m and is on a slope θ above the horizontal. At a distance x down the slope below him is a ledge. He is tied in to a climbing rope and being belayed from above, so that if he slips he won't simply plunge to his death. Climbing ropes are intentionally made out of stretchy material so that in a fall, the climber gets a gentle catch rather than a violent force that would hurt. However, the rope should not be more stretchy than necessary because of situations like this one: if the rope were to stretch by more than x , the climber would hit the ledge.

- (a) Find the spring constant that the rope should have in order to limit the amount of rope stretch to x . ✓
- (b) Show that your answer to part a has the right units.
- (c) Analyze the mathematical dependence of the result on each of the variables, and verify that it makes sense physically.



Problem 91.

Key to symbols:

■ easy ■ typical ■ challenging ■ difficult ■ very difficult

✓ An answer check is available at www.lightandmatter.com.

Exercises

Exercise 3A: Force and Motion

Equipment:

2-meter pieces of butcher paper

wood blocks with hooks

string

masses to put on top of the blocks to increase friction

spring scales (preferably calibrated in Newtons)

Suppose a person pushes a crate, sliding it across the floor at a certain speed, and then repeats the same thing but at a higher speed. This is essentially the situation you will act out in this exercise. What do you think is different about her force on the crate in the two situations? Discuss this with your group and write down your hypothesis:

1. First you will measure the amount of friction between the wood block and the butcher paper when the wood and paper surfaces are slipping over each other. The idea is to attach a spring scale to the block and then slide the butcher paper under the block while using the scale to keep the block from moving with it. Depending on the amount of force your spring scale was designed to measure, you may need to put an extra mass on top of the block in order to increase the amount of friction. It is a good idea to use long piece of string to attach the block to the spring scale, since otherwise one tends to pull at an angle instead of directly horizontally.

First measure the amount of friction force when sliding the butcher paper as slowly as possible:

Now measure the amount of friction force at a significantly higher speed, say 1 meter per second. (If you try to go too fast, the motion is jerky, and it is impossible to get an accurate reading.)

Discuss your results. Why are we justified in assuming that the string's force on the block (i.e., the scale reading) is the same amount as the paper's frictional force on the block?

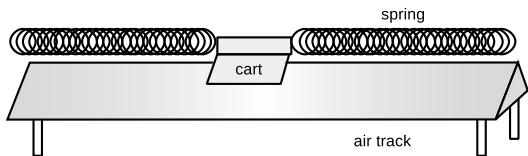
2. Now try the same thing but with the block moving and the paper standing still. Try two different speeds.

Do your results agree with your original hypothesis? If not, discuss what's going on. How does the block "know" how fast to go?

Exercise 3B: Vibrations

Equipment:

- air track and carts of two different masses
- springs
- spring scales

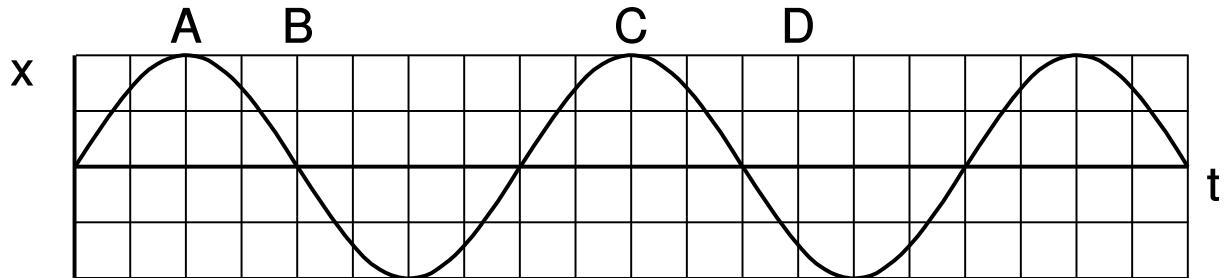


Place the cart on the air track and attach springs so that it can vibrate.

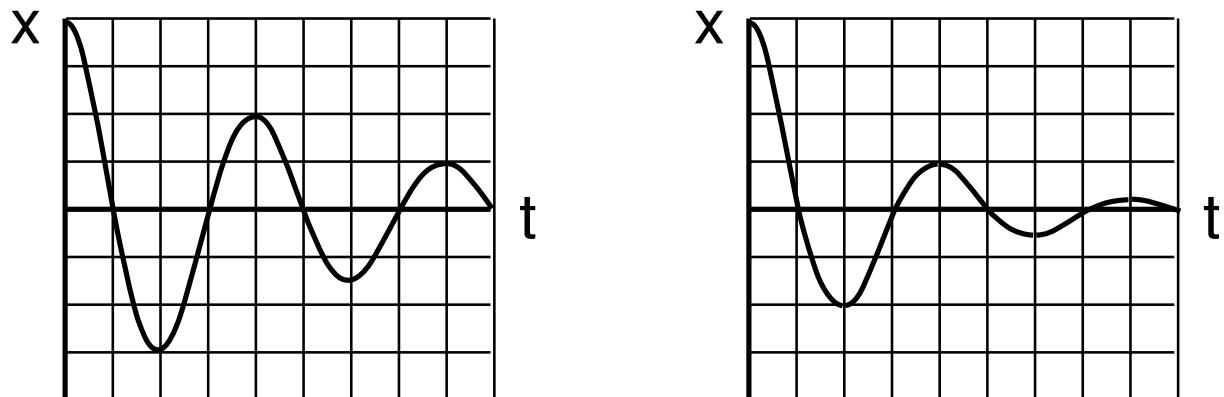
1. Test whether the period of vibration depends on amplitude. Try at least two moderate amplitudes, for which the springs do not go slack, and at least one amplitude that is large enough so that they do go slack.
2. Try a cart with a different mass. Does the period change by the expected factor, based on the equation $T = 2\pi\sqrt{m/k}$?
3. Use a spring scale to pull the cart away from equilibrium, and make a graph of force versus position. Is it linear? If so, what is its slope?
4. Test the equation $T = 2\pi\sqrt{m/k}$ numerically.

Exercise 3C: Worksheet on Resonance

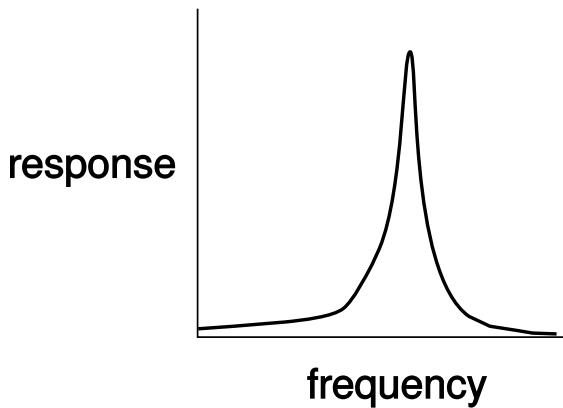
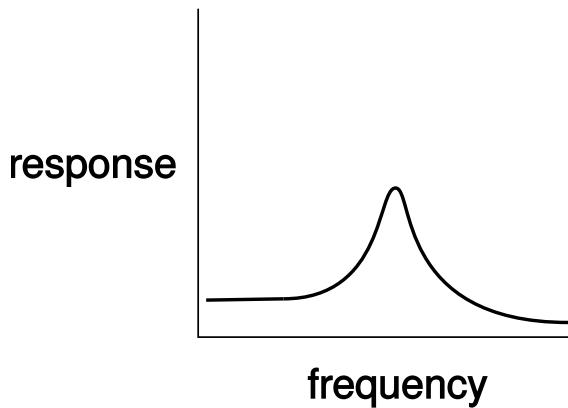
1. Compare the oscillator's energies at A, B, C, and D.



2. Compare the Q values of the two oscillators.



3. Match the x-t graphs in #2 with the amplitude-frequency graphs below.



Exercise D is on the following two pages.

Exercise 3D: Vectors and Motion

Each diagram on page 249 shows the motion of an object in an $x - y$ plane. Each dot is one location of the object at one moment in time. The time interval from one dot to the next is always the same, so you can think of the vector that connects one dot to the next as a \mathbf{v} vector, and subtract to find $\Delta\mathbf{v}$ vectors.

1. Suppose the object in diagram 1 is moving from the top left to the bottom right. Deduce whatever you can about the force acting on it. Does the force always have the same magnitude? The same direction?

Invent a physical situation that this diagram could represent.

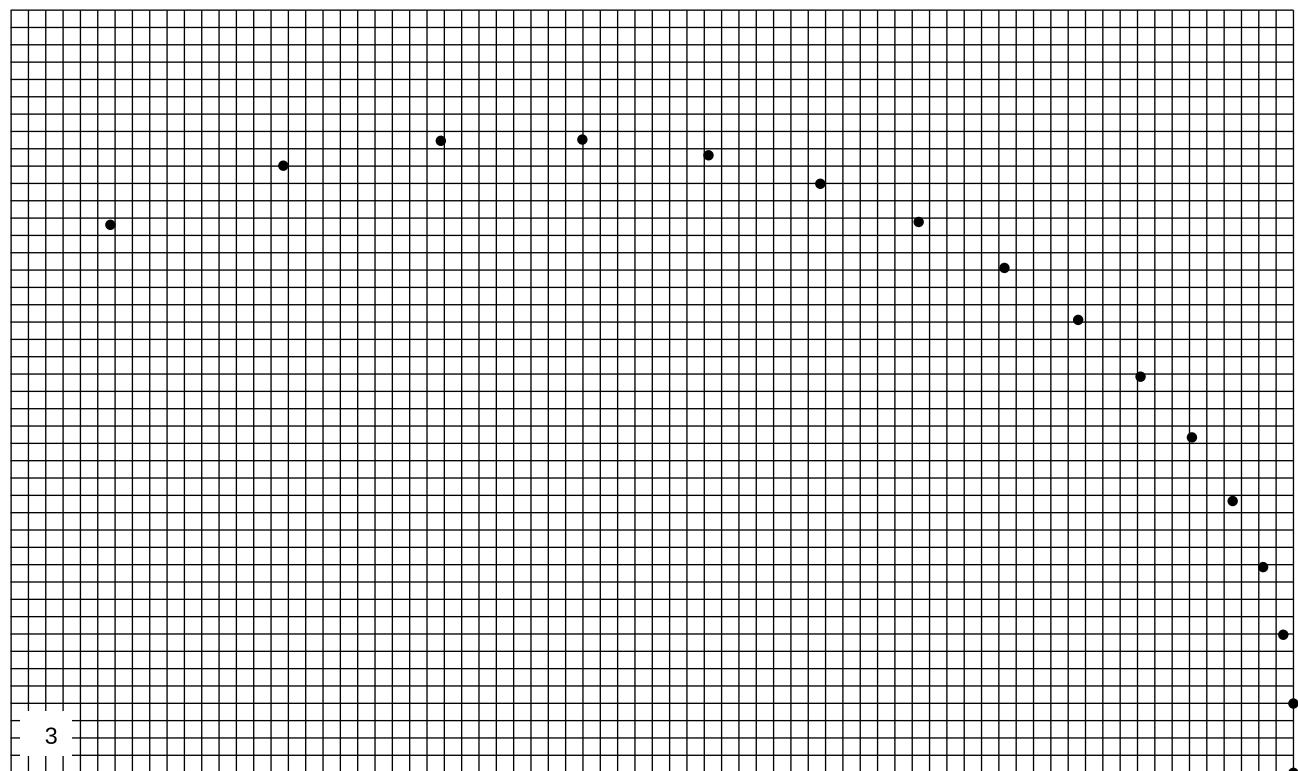
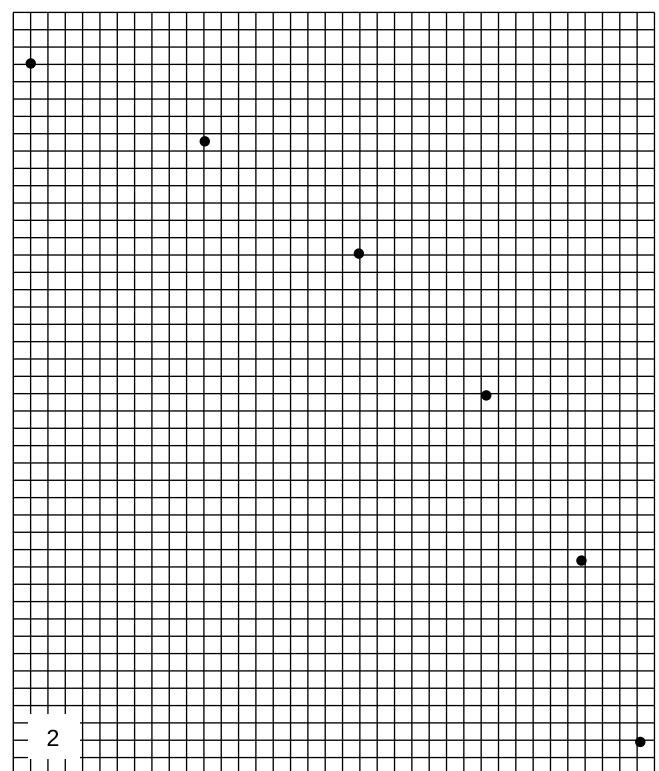
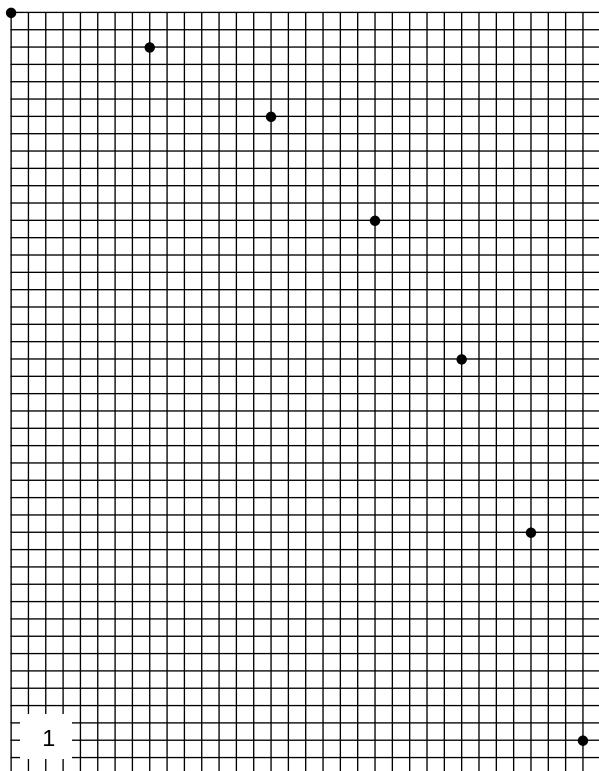
What if you reinterpret the diagram, and reverse the object's direction of motion?

2. What can you deduce about the force that is acting in diagram 2?

Invent a physical situation that diagram 2 could represent.

3. What can you deduce about the force that is acting in diagram 3?

Invent a physical situation.



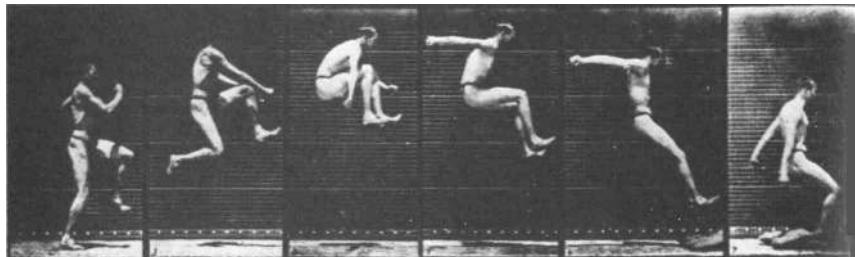
Chapter 4

Conservation of Angular Momentum

4.1 Angular momentum in two dimensions

4.1.1 Angular momentum

“Sure, and maybe the sun won’t come up tomorrow.” Of course, the sun only appears to go up and down because the earth spins, so the cliche should really refer to the unlikelihood of the earth’s stopping its rotation abruptly during the night. Why can’t it stop? It wouldn’t violate conservation of momentum, because the earth’s rotation doesn’t add anything to its momentum. While California spins in one direction, some equally massive part of India goes the opposite way, canceling its momentum. A halt to Earth’s rotation would entail a drop in kinetic energy, but that energy could simply be converted into some other form, such as heat.



a / The jumper can’t move his legs counterclockwise without moving his arms clockwise.
(Thomas Eakins.)

Other examples along these lines are not hard to find. An atom spins at the same rate for billions of years. A high-diver who is rotating when he comes off the board does not need to make any physical effort to continue rotating, and indeed would be unable to stop rotating before he hit the water.

These observations have the hallmarks of a conservation law:

A closed system is involved. Nothing is making an effort to twist the earth, the hydrogen atom, or the high-diver. They are isolated from rotation-changing influences, i.e., they are closed systems.

Something remains unchanged. There appears to be a numerical quantity for measuring rotational motion such that the total amount of that quantity remains constant in a closed system.

Something can be transferred back and forth without changing the total amount. In the photo of the old-fashioned high jump, a, the jumper wants to get his feet out in front of him so he can keep from doing a “face plant” when he lands. Bringing his feet forward would involve a certain quantity of counterclockwise rotation, but he didn’t start out with any rotation when he left the ground. Suppose we consider counterclockwise as positive and clockwise as negative. The only way his legs can acquire some positive rotation is if some other part of his body picks up an equal amount of negative rotation. This is why he swings his arms up behind him, clockwise.

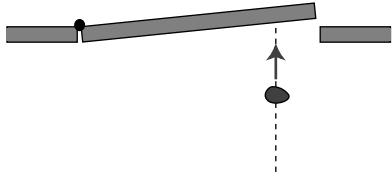
What numerical measure of rotational motion is conserved? Car engines and old-fashioned LP records have speeds of rotation measured in rotations per minute (r.p.m.), but the number of rotations per minute (or per second) is not a conserved quantity. A twirling figure skater, for instance, can pull her arms in to increase her r.p.m.’s. The first section of this chapter deals with the numerical definition of the quantity of rotation that results in a valid conservation law.

When most people think of rotation, they think of a solid object like a wheel rotating in a circle around a fixed point. Examples of this type of rotation, called rigid rotation or rigid-body rotation, include a spinning top, a seated child’s swinging leg, and a helicopter’s spinning propeller. Rotation, however, is a much more general phenomenon, and includes noncircular examples such as a comet in an elliptical orbit around the sun, or a cyclone, in which the core completes a circle more quickly than the outer parts.

If there is a numerical measure of rotational motion that is a conserved quantity, then it must include nonrigid cases like these, since nonrigid rotation can be traded back and forth with rigid rotation. For instance, there is a trick for finding out if an egg is raw or hardboiled. If you spin a hardboiled egg and then stop it briefly with your finger, it stops dead. But if you do the same with a raw egg, it springs back into rotation because the soft interior was still swirling around within the momentarily motionless shell. The pattern of flow of the liquid part is presumably very complex and nonuniform due to the asymmetric shape of the egg and the different consistencies of the yolk and the white, but there is apparently some way to describe the liquid’s total amount of rotation with a single number, of which some percentage is given back to the shell when you release it.

The best strategy is to devise a way of defining the amount of rotation of a single small part of a system. The amount of rotation of a system such as a cyclone will then be defined as the total of all the contributions from its many small parts.

The quest for a conserved quantity of rotation even requires us to broaden the rotation concept to include cases where the motion



b / An overhead view of a piece of putty being thrown at a door. Even though the putty is neither spinning nor traveling along a curve, we must define it has having some kind of “rotation” because it is able to make the door rotate.

doesn't repeat or even curve around. If you throw a piece of putty at a door, b, the door will recoil and start rotating. The putty was traveling straight, not in a circle, but if there is to be a general conservation law that can cover this situation, it appears that we must describe the putty as having had some "rotation," which it then gave up to the door. The best way of thinking about it is to attribute rotation to any moving object or part of an object that changes its angle in relation to the axis of rotation. In the putty-and-door example, the hinge of the door is the natural point to think of as an axis, and the putty changes its angle as seen by someone standing at the hinge, c. For this reason, the conserved quantity we are investigating is called *angular momentum*. The symbol for angular momentum can't be "a" or "m," since those are used for acceleration and mass, so the letter *L* is arbitrarily chosen instead.

Imagine a 1 kg blob of putty, thrown at the door at a speed of 1 m/s, which hits the door at a distance of 1 m from the hinge. We define this blob to have 1 unit of angular momentum. When it hits the door, the door will recoil and start rotating. We can use the speed at which the door recoils as a measure of the angular momentum the blob brought in.¹

Experiments show, not surprisingly, that a 2 kg blob thrown in the same way makes the door rotate twice as fast, so the angular momentum of the putty blob must be proportional to mass,

$$L \propto m.$$

Similarly, experiments show that doubling the velocity of the blob will have a doubling effect on the result, so its angular momentum must be proportional to its velocity as well,

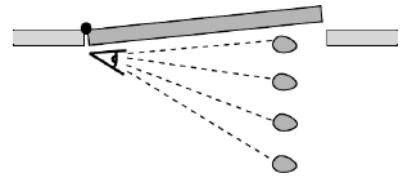
$$L \propto mv.$$

You have undoubtedly had the experience of approaching a closed door with one of those bar-shaped handles on it and pushing on the wrong side, the side close to the hinges. You feel like an idiot, because you have so little leverage that you can hardly budge the door. The same would be true with the putty blob. Experiments would show that the amount of rotation the blob can give to the door is proportional to the distance, *r*, from the axis of rotation, so angular momentum must be proportional to *r* as well,

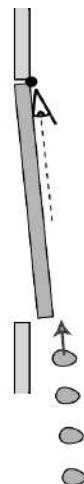
$$L \propto mvr.$$

We are almost done, but there is one missing ingredient. We know on grounds of symmetry that a putty ball thrown directly

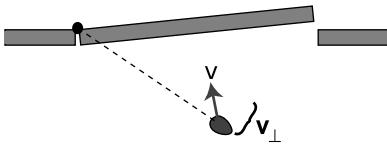
¹We assume that the door is much more massive than the blob. Under this assumption, the speed at which the door recoils is much less than the original speed of the blob, so the blob has lost essentially all its angular momentum, and given it to the door.



c / As seen by someone standing at the axis, the putty changes its angular position. We therefore define it as having angular momentum.



d / A putty blob thrown directly at the axis has no angular motion, and therefore no angular momentum. It will not cause the door to rotate.



e / Only the component of the velocity vector perpendicular to the line connecting the object to the axis should be counted into the definition of angular momentum.



f / A figure skater pulls in her arms so that she can execute a spin more rapidly.

inward toward the hinge will have no angular momentum to give to the door. After all, there would not even be any way to decide whether the ball's rotation was clockwise or counterclockwise in this situation. It is therefore only the component of the blob's velocity vector perpendicular to the door that should be counted in its angular momentum,

$$L = mv_{\perp}r.$$

More generally, v_{\perp} should be thought of as the component of the object's velocity vector that is perpendicular to the line joining the object to the axis of rotation.

We find that this equation agrees with the definition of the original putty blob as having one unit of angular momentum, and we can now see that the units of angular momentum are $(\text{kg}\cdot\text{m}/\text{s})\cdot\text{m}$, i.e., $\text{kg}\cdot\text{m}^2/\text{s}$. Summarizing, we have

$L = mv_{\perp}r$ [angular momentum of a particle in two dimensions], where m is the particle's mass, v_{\perp} is the component of its velocity vector perpendicular to the line joining it to the axis of rotation, and r is its distance from the axis. (Note that r is not necessarily the radius of a circle.) Positive and negative signs of angular momentum are used to describe opposite directions of rotation. The angular momentum of a finite-sized object or a system of many objects is found by dividing it up into many small parts, applying the equation to each part, and adding to find the total amount of angular momentum. (As implied by the word "particle," matter isn't the only thing that can have angular momentum. Light can also have angular momentum, and the above equation would not apply to light.)

Conservation of angular momentum has been verified over and over again by experiment, and is now believed to be one of the most fundamental principles of physics, along with conservation of mass, energy, and momentum.

A figure skater pulls her arms in.

example 1

When a figure skater is twirling, there is very little friction between her and the ice, so she is essentially a closed system, and her angular momentum is conserved. If she pulls her arms in, she is decreasing r for all the atoms in her arms. It would violate conservation of angular momentum if she then continued rotating at the same speed, i.e., taking the same amount of time for each revolution, because her arms' contributions to her angular momentum would have decreased, and no other part of her would have increased its angular momentum. This is impossible because it would violate conservation of angular momentum. If her total angular momentum is to remain constant, the decrease in r for her arms must be compensated for by an overall increase in her rate of rotation. That is, by pulling her arms in, she substantially reduces the time for each rotation.

Earth's slowing rotation and the receding moon example 2

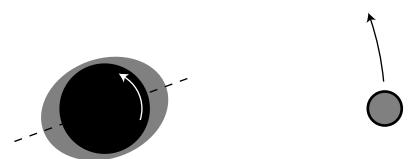
The earth's rotation is actually slowing down very gradually, with the kinetic energy being dissipated as heat by friction between the land and the tidal bulges raised in the seas by the earth's gravity. Does this mean that angular momentum is not really perfectly conserved? No, it just means that the earth is not quite a closed system by itself. If we consider the earth and moon as a system, then the angular momentum lost by the earth must be gained by the moon somehow. In fact very precise measurements of the distance between the earth and the moon have been carried out by bouncing laser beams off of a mirror left there by astronauts, and these measurements show that the moon is receding from the earth at a rate of 4 centimeters per year! The moon's greater value of r means that it has a greater angular momentum, and the increase turns out to be exactly the amount lost by the earth. In the days of the dinosaurs, the days were significantly shorter, and the moon was closer and appeared bigger in the sky.

But what force is causing the moon to speed up, drawing it out into a larger orbit? It is the gravitational forces of the earth's tidal bulges. In figure g, the earth's rotation is counterclockwise (arrow). The moon's gravity creates a bulge on the side near it, because its gravitational pull is stronger there, and an "anti-bulge" on the far side, since its gravity there is weaker. For simplicity, let's focus on the tidal bulge closer to the moon. Its frictional force is trying to slow down the earth's rotation, so its force on the earth's solid crust is toward the bottom of the figure. By Newton's third law, the crust must thus make a force on the bulge which is toward the top of the figure. This causes the bulge to be pulled forward at a slight angle, and the bulge's gravity therefore pulls the moon forward, accelerating its orbital motion about the earth and flinging it outward.

The result would obviously be extremely difficult to calculate directly, and this is one of those situations where a conservation law allows us to make precise quantitative statements about the outcome of a process when the calculation of the process itself would be prohibitively complex.

Restriction to rotation in a plane

Is angular momentum a vector, or a scalar? It does have a direction in space, but it's a direction of rotation, not a straight-line direction like the directions of vectors such as velocity or force. It turns out that there is a way of defining angular momentum as a vector, but in this section the examples will be confined to a single plane of rotation, i.e., effectively two-dimensional situations. In this special case, we can choose to visualize the plane of rotation from one side or the other, and to define clockwise and counterclockwise rotation as having opposite signs of angular momentum. "Efec-



g / A view of the earth-moon system from above the north pole. All distances have been highly distorted for legibility.

tively” two-dimensional means that we can deal with objects that aren’t flat, as long as the velocity vectors of all their parts lie in a plane.

Discussion Questions

A Conservation of plain old momentum, p , can be thought of as the greatly expanded and modified descendant of Galileo’s original principle of inertia, that no force is required to keep an object in motion. The principle of inertia is counterintuitive, and there are many situations in which it appears superficially that a force *is* needed to maintain motion, as maintained by Aristotle. Think of a situation in which conservation of angular momentum, L , also seems to be violated, making it seem incorrectly that something external must act on a closed system to keep its angular momentum from “running down.”

4.1.2 Application to planetary motion

We now discuss the application of conservation of angular momentum to planetary motion, both because of its intrinsic importance and because it is a good way to develop a visual intuition for angular momentum.

Kepler’s law of equal areas states that the area swept out by a planet in a certain length of time is always the same. Angular momentum had not been invented in Kepler’s time, and he did not even know the most basic physical facts about the forces at work. He thought of this law as an entirely empirical and unexpectedly simple way of summarizing his data, a rule that succeeded in describing and predicting how the planets sped up and slowed down in their elliptical paths. It is now fairly simple, however, to show that the equal area law amounts to a statement that the planet’s angular momentum stays constant.

There is no simple geometrical rule for the area of a pie wedge cut out of an ellipse, but if we consider a very short time interval, as shown in figure h, the shaded shape swept out by the planet is very nearly a triangle. We do know how to compute the area of a triangle. It is one half the product of the base and the height:

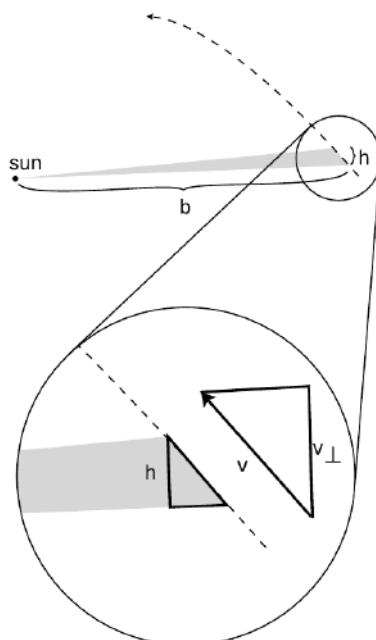
$$\text{area} = \frac{1}{2}bh.$$

We wish to relate this to angular momentum, which contains the variables r and v_{\perp} . If we consider the sun to be the axis of rotation, then the variable r is identical to the base of the triangle, $r = b$. Referring to the magnified portion of the figure, v_{\perp} can be related to h , because the two right triangles are similar:

$$\frac{h}{\text{distance traveled}} = \frac{v_{\perp}}{|\mathbf{v}|}$$

The area can thus be rewritten as

$$\text{area} = \frac{1}{2}r \frac{v_{\perp}(\text{distance traveled})}{|\mathbf{v}|}.$$



h / The area swept out by a planet in its orbit.

The distance traveled equals $|\mathbf{v}|\Delta t$, so this simplifies to

$$\text{area} = \frac{1}{2}rv_{\perp}\Delta t.$$

We have found the following relationship between angular momentum and the rate at which area is swept out:

$$L = 2m \frac{\text{area}}{\Delta t}.$$

The factor of 2 in front is simply a matter of convention, since any conserved quantity would be an equally valid conserved quantity if you multiplied it by a constant. The factor of m was not relevant to Kepler, who did not know the planets' masses, and who was only describing the motion of one planet at a time.

We thus find that Kepler's equal-area law is equivalent to a statement that the planet's angular momentum remains constant. But wait, why should it remain constant? — the planet is not a closed system, since it is being acted on by the sun's gravitational force. There are two valid answers. The first is that it is actually the total angular momentum of the sun plus the planet that is conserved. The sun, however, is millions of times more massive than the typical planet, so it accelerates very little in response to the planet's gravitational force. It is thus a good approximation to say that the sun doesn't move at all, so that no angular momentum is transferred between it and the planet.

The second answer is that to change the planet's angular momentum requires not just a force but a force applied in a certain way. Later in this section (starting on page 260) we discuss the transfer of angular momentum by a force, but the basic idea here is that a force directly in toward the axis does not change the angular momentum.

Discussion Questions

A Suppose an object is simply traveling in a straight line at constant speed. If we pick some point not on the line and call it the axis of rotation, is area swept out by the object at a constant rate?

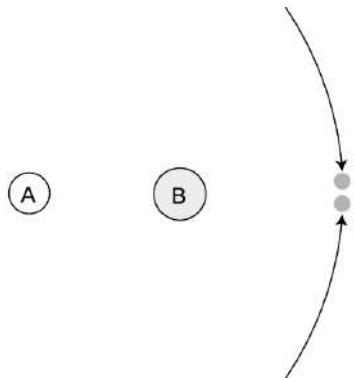
B The figure is a strobe photo of a pendulum bob, taken from underneath the pendulum looking straight up. The black string can't be seen in the photograph. The bob was given a slight sideways push when it was released, so it did not swing in a plane. The bright spot marks the center, i.e., the position the bob would have if it hung straight down at us. Does the bob's angular momentum appear to remain constant if we consider the center to be the axis of rotation?

4.1.3 Two theorems about angular momentum

With plain old momentum, \mathbf{p} , we had the freedom to work in any inertial frame of reference we liked. The same object could have different values of momentum in two different frames, if the



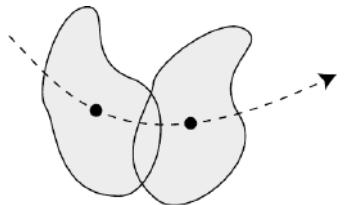
i / Discussion question B.



j / Two asteroids collide.



k / Everyone has a strong tendency to think of the diver as rotating about his own center of mass. However, he is flying in an arc, and he also has angular momentum because of this motion.



l / This rigid object has angular momentum both because it is spinning about its center of mass and because it is moving through space.

frames were not at rest with respect to each other. Conservation of momentum, however, would be true in either frame. As long as we employed a single frame consistently throughout a calculation, everything would work.

The same is true for angular momentum, and in addition there is an ambiguity that arises from the definition of an axis of rotation. For a wheel, the natural choice of an axis of rotation is obviously the axle, but what about an egg rotating on its side? The egg has an asymmetric shape, and thus no clearly defined geometric center. A similar issue arises for a cyclone, which does not even have a sharply defined shape, or for a complicated machine with many gears. The following theorem, the first of two presented in this section, explains how to deal with this issue. Although I have put descriptive titles above both theorems, they have no generally accepted names. The proofs, given on page 1028, use the vector cross-product technique introduced in section 4.3, which greatly simplifies them.

The choice of axis theorem: It is entirely arbitrary what point one defines as the axis for purposes of calculating angular momentum. If a closed system's angular momentum is conserved when calculated with one choice of axis, then it will be conserved for any other choice of axis. Likewise, any inertial frame of reference may be used. The theorem also holds in the case where the system is not closed, but the total external force is zero.

Colliding asteroids described with different axes example 3

Observers on planets A and B both see the two asteroids colliding. The asteroids are of equal mass and their impact speeds are the same. Astronomers on each planet decide to define their own planet as the axis of rotation. Planet A is twice as far from the collision as planet B. The asteroids collide and stick. For simplicity, assume planets A and B are both at rest.

With planet A as the axis, the two asteroids have the same amount of angular momentum, but one has positive angular momentum and the other has negative. Before the collision, the total angular momentum is therefore zero. After the collision, the two asteroids will have stopped moving, and again the total angular momentum is zero. The total angular momentum both before and after the collision is zero, so angular momentum is conserved if you choose planet A as the axis.

The only difference with planet B as axis is that r is smaller by a factor of two, so all the angular momenta are halved. Even though the angular momenta are different than the ones calculated by planet A, angular momentum is still conserved.

The earth spins on its own axis once a day, but simultaneously travels in its circular one-year orbit around the sun, so any given part of it traces out a complicated loopy path. It would seem difficult

to calculate the earth's angular momentum, but it turns out that there is an intuitively appealing shortcut: we can simply add up the angular momentum due to its spin plus that arising from its center of mass's circular motion around the sun. This is a special case of the following general theorem:

The spin theorem: An object's angular momentum with respect to some outside axis A can be found by adding up two parts:

- (1) The first part is the object's angular momentum found by using its own center of mass as the axis, i.e., the angular momentum the object has because it is spinning.
- (2) The other part equals the angular momentum that the object would have with respect to the axis A if it had all its mass concentrated at and moving with its center of mass.

A system with its center of mass at rest *example 4*

In the special case of an object whose center of mass is at rest, the spin theorem implies that the object's angular momentum is the same regardless of what axis we choose. (This is an even stronger statement than the choice of axis theorem, which only guarantees that angular momentum is conserved for any given choice of axis, without specifying that it is the same for all such choices.)

Angular momentum of a rigid object *example 5*

▷ A motorcycle wheel has almost all its mass concentrated at the outside. If the wheel has mass m and radius r , and the time required for one revolution is T , what is the spin part of its angular momentum?

▷ This is an example of the commonly encountered special case of rigid motion, as opposed to the rotation of a system like a hurricane in which the different parts take different amounts of time to go around. We don't really have to go through a laborious process of adding up contributions from all the many parts of a wheel, because they are all at about the same distance from the axis, and are all moving around the axis at about the same speed. The velocity is all perpendicular to the spokes,

$$\begin{aligned} v_{\perp} &= (\text{circumference})/T \\ &= 2\pi r/T \end{aligned}$$

and the angular momentum of the wheel about its center is

$$\begin{aligned} L &= mv_{\perp}r \\ &= m(2\pi r/T)r \\ &= 2\pi mr^2/T. \end{aligned}$$

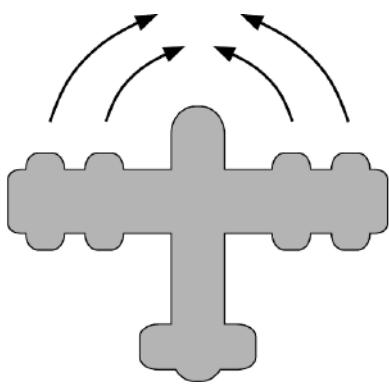
Note that although the factors of 2π in this expression is peculiar to a wheel with its mass concentrated on the rim, the proportionality to m/T would have been the same for any other rigidly rotating

object. Although an object with a noncircular shape does not have a radius, it is also true in general that angular momentum is proportional to the square of the object's size for fixed values of m and T . For instance doubling an object's size doubles both the v_{\perp} and r factors in the contribution of each of its parts to the total angular momentum, resulting in an overall factor of four increase.

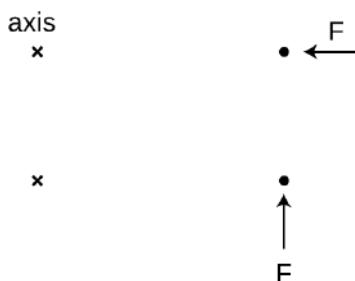
4.1.4 Torque

Force is the rate of transfer of momentum. The corresponding quantity in the case of angular momentum is called torque (rhymes with "fork"). Where force tells us how hard we are pushing or pulling on something, torque indicates how hard we are twisting on it. Torque is represented by the Greek letter tau, τ , and the rate of change of an object's angular momentum equals the total torque acting on it,

$$\tau_{total} = dL/dt.$$



m / The plane's four engines produce zero total torque but not zero total force.



n / The simple physical situation we use to derive an equation for torque. A force that points directly in at or out away from the axis produces neither clockwise nor counterclockwise angular momentum. A force in the perpendicular direction does transfer angular momentum.

As with force and momentum, it often happens that angular momentum recedes into the background and we focus our interest on the torques. The torque-focused point of view is exemplified by the fact that many scientifically untrained but mechanically apt people know all about torque, but none of them have heard of angular momentum. Car enthusiasts eagerly compare engines' torques, and there is a tool called a torque wrench which allows one to apply a desired amount of torque to a screw and avoid overtightening it.

Torque distinguished from force

Of course a force is necessary in order to create a torque — you can't twist a screw without pushing on the wrench — but force and torque are two different things. One distinction between them is direction. We use positive and negative signs to represent forces in the two possible directions along a line. The direction of a torque, however, is clockwise or counterclockwise, not a linear direction.

The other difference between torque and force is a matter of leverage. A given force applied at a door's knob will change the door's angular momentum twice as rapidly as the same force applied halfway between the knob and the hinge. The same amount of force produces different amounts of torque in these two cases.

It's possible to have a zero total torque with a nonzero total force. An airplane with four jet engines would be designed so that their forces are balanced on the left and right. Their forces are all in the same direction, but the clockwise torques of two of the engines are canceled by the counterclockwise torques of the other two, giving zero total torque.

Conversely we can have zero total force and nonzero total torque. A merry-go-round's engine needs to supply a nonzero torque on it to bring it up to speed, but there is zero total force on it. If there

was not zero total force on it, its center of mass would accelerate!

Relationship between force and torque

How do we calculate the amount of torque produced by a given force? Since it depends on leverage, we should expect it to depend on the distance between the axis and the point of application of the force. I'll work out an equation relating torque to force for a particular very simple situation, and give a more rigorous derivation on page 290, after developing some mathematical techniques that dramatically shorten and simplify the proof.

Consider a pointlike object which is initially at rest at a distance r from the axis we have chosen for defining angular momentum. We first observe that a force directly inward or outward, along the line connecting the axis to the object, does not impart any angular momentum to the object.

A force perpendicular to the line connecting the axis and the object does, however, make the object pick up angular momentum. Newton's second law gives

$$a = F/m,$$

and using $a = dv/dt$ we find the velocity the object acquires after a time dt ,

$$dv = F dt / m.$$

We're trying to relate force to a change in angular momentum, so we multiply both sides of the equation by mr to give

$$\begin{aligned} m dv r &= F dt r \\ dL &= F dt r. \end{aligned}$$

Dividing by dt gives the torque:

$$\begin{aligned} \frac{dL}{dt} &= Fr \\ \tau &= Fr. \end{aligned}$$

If a force acts at an angle other than 0 or 90° with respect to the line joining the object and the axis, it would be only the component of the force perpendicular to the line that would produce a torque,

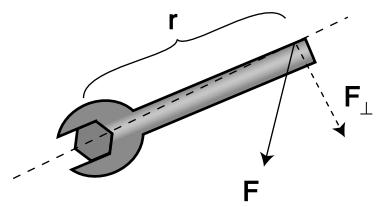
$$\tau = F_{\perp} r.$$

Although this result was proved under a simplified set of circumstances, it is more generally valid:²

Relationship between force and torque: The rate at which a force transfers angular momentum to an object, i.e., the torque produced by the force, is given by

$$|\tau| = r |F_{\perp}|,$$

²A proof is given in example 28 on page 290



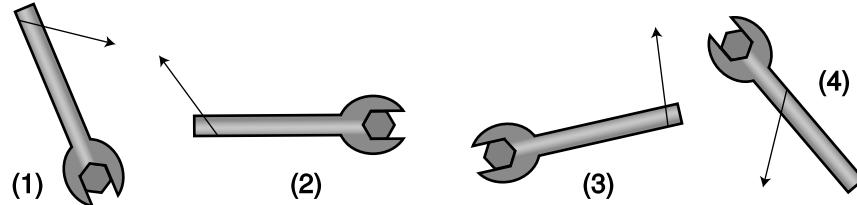
o / The geometric relationships referred to in the relationship between force and torque.

where r is the distance from the axis to the point of application of the force, and F_{\perp} is the component of the force that is perpendicular to the line joining the axis to the point of application.

The equation is stated with absolute value signs because the positive and negative signs of force and torque indicate different things, so there is no useful relationship between them. The sign of the torque must be found by physical inspection of the case at hand.

From the equation, we see that the units of torque can be written as newtons multiplied by meters. Metric torque wrenches are calibrated in N·m, but American ones use foot-pounds, which is also a unit of distance multiplied by a unit of force. We know from our study of mechanical work that newtons multiplied by meters equal joules, but torque is a completely different quantity from work, and nobody writes torques with units of joules, even though it would be technically correct.

p / Self-check.



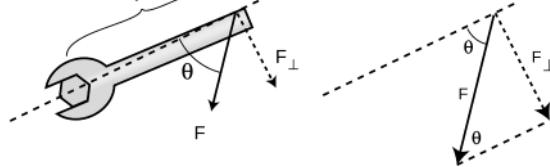
self-check A

Compare the magnitudes and signs of the four torques shown in figure p.

▷ Answer, p. 1060

How torque depends on the direction of the force example 6

- ▷ How can the torque applied to the wrench in the figure be expressed in terms of r , $|F|$, and the angle θ ?
- ▷ The force vector and its F_{\perp} component form the hypotenuse and one leg of a right triangle,



and the interior angle opposite to F_{\perp} equals θ . The absolute value of F_{\perp} can thus be expressed as

$$F_{\perp} = |F| \sin \theta,$$

leading to

$$|\tau| = r|F| \sin \theta.$$

Sometimes torque can be more neatly visualized in terms of the quantity r_{\perp} shown in the figure on the left, which gives us a third way of expressing the relationship between torque and force:

$$|\tau| = r_{\perp} |F|.$$

Of course you wouldn't want to go and memorize all three equations for torque. Starting from any one of them you could easily derive the other two using trigonometry. Familiarizing yourself with them can however clue you in to easier avenues of attack on certain problems.

The torque due to gravity

Up until now we've been thinking in terms of a force that acts at a single point on an object, such as the force of your hand on the wrench. This is of course an approximation, and for an extremely realistic calculation of your hand's torque on the wrench you might need to add up the torques exerted by each square millimeter where your skin touches the wrench. This is seldom necessary. But in the case of a gravitational force, there is never any single point at which the force is applied. Our planet is exerting a separate tug on every brick in the Leaning Tower of Pisa, and the total gravitational torque on the tower is the sum of the torques contributed by all the little forces. Luckily there is a trick that allows us to avoid such a massive calculation. It turns out that for purposes of computing the total gravitational torque on an object, you can get the right answer by just pretending that the whole gravitational force acts at the object's center of mass.

Gravitational torque on an outstretched arm *example 7*

▷ Your arm has a mass of 3.0 kg, and its center of mass is 30 cm from your shoulder. What is the gravitational torque on your arm when it is stretched out horizontally to one side, taking the shoulder to be the axis?

▷ The total gravitational force acting on your arm is

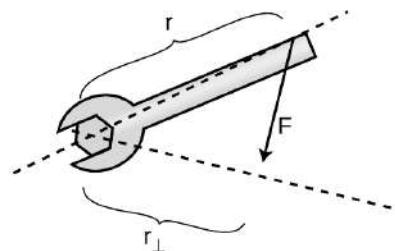
$$|\mathbf{F}| = (3.0 \text{ kg})(9.8 \text{ m/s}^2) = 29 \text{ N}.$$

For the purpose of calculating the gravitational torque, we can treat the force as if it acted at the arm's center of mass. The force is straight down, which is perpendicular to the line connecting the shoulder to the center of mass, so

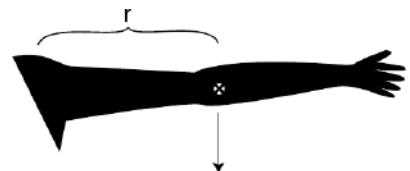
$$F_{\perp} = |\mathbf{F}| = 29 \text{ N}.$$

Continuing to pretend that the force acts at the center of the arm, r equals $30 \text{ cm} = 0.30 \text{ m}$, so the torque is

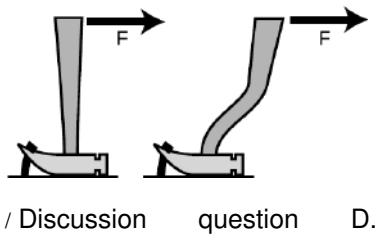
$$\tau = r F_{\perp} = 9 \text{ N} \cdot \text{m}.$$



q / Visualizing torque in terms of r_{\perp} .



r / Example 7.



Discussion Questions

A This series of discussion questions deals with past students' incorrect reasoning about the following problem.

Suppose a comet is at the point in its orbit shown in the figure. The only force on the comet is the sun's gravitational force. Throughout the question, define all torques and angular momenta using the sun as the axis.

- (1) *Is the sun producing a nonzero torque on the comet? Explain.*
- (2) *Is the comet's angular momentum increasing, decreasing, or staying the same? Explain.*

Explain what is wrong with the following answers. In some cases, the answer is correct, but the reasoning leading up to it is wrong.

- (a) Incorrect answer to part (1): "Yes, because the sun is exerting a force on the comet, and the comet is a certain distance from the sun."
- (b) Incorrect answer to part (1): "No, because the torques cancel out."
- (c) Incorrect answer to part (2): "Increasing, because the comet is speeding up."



u / Discussion question A.

B You whirl a rock over your head on the end of a string, and gradually pull in the string, eventually cutting the radius in half. What happens to the rock's angular momentum? What changes occur in its speed, the time required for one revolution, and its acceleration? Why might the string break?

C A helicopter has, in addition to the huge fan blades on top, a smaller propeller mounted on the tail that rotates in a vertical plane. Why?

D Which claw hammer would make it easier to get the nail out of the wood if the same force was applied in the same direction?

E The photo shows an amusement park ride whose two cars rotate in opposite directions. Why is this a good design?

4.1.5 Applications to statics

In chapter 2 I defined equilibrium as a situation where the interaction energy is minimized. This is the same as a condition of zero total force, or constant momentum. Thus a car is in equilibrium not just when it is parked but also when it is cruising down a straight road with constant momentum.

Likewise there are many cases where a system is not closed but maintains constant angular momentum. When a merry-go-round is running at constant angular momentum, the engine's torque is

being canceled by the torque due to friction.

It's not enough for a boat not to sink — we'd also like to avoid having it capsize. For this reason, we now redefine equilibrium as follows.

When an object has constant momentum and constant angular momentum, we say that it is in equilibrium. Again, this is a scientific redefinition of the common English word, since in ordinary speech nobody would describe a car spinning out on an icy road as being in equilibrium.

Very commonly, however, we are interested in cases where an object is not only in equilibrium but also at rest, and this corresponds more closely to the usual meaning of the word. Statics is the branch of physics concerned with problems such as these.

Solving statics problems is now simply a matter of applying and combining some things you already know:

- You know the behaviors of the various types of forces, for example that a frictional force is always parallel to the surface of contact.
- You know about vector addition of forces. It is the vector sum of the forces that must equal zero to produce equilibrium.
- You know about torque. The total torque acting on an object must be zero if it is to be in equilibrium.
- You know that the choice of axis is arbitrary, so you can make a choice of axis that makes the problem easy to solve.

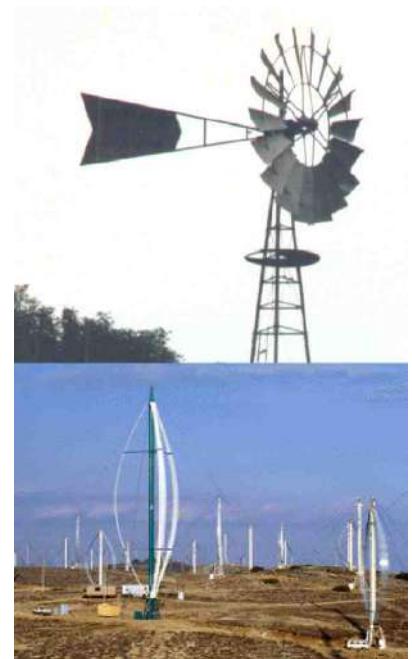
In general, this type of problem could involve four equations in four unknowns: three equations that say the force components add up to zero, and one equation that says the total torque is zero. Most cases you'll encounter will not be this complicated. In the example below, only the equation for zero total torque is required in order to get an answer.

A flagpole

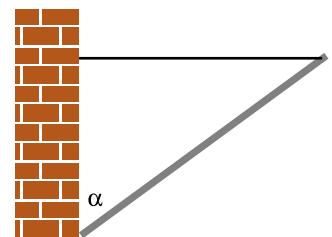
example 8

▷ A 10-kg flagpole is being held up by a lightweight horizontal cable, and is propped against the foot of a wall as shown in the figure. If the cable is only capable of supporting a tension of 70 N, how great can the angle α be without breaking the cable?

▷ All three objects in the figure are supposed to be in equilibrium: the pole, the cable, and the wall. Whichever of the three objects we pick to investigate, all the forces and torques on it have to cancel out. It is not particularly helpful to analyze the forces and torques on the wall, since it has forces on it from the ground that are not given and that we don't want to find. We could study the forces and torques on the cable, but that doesn't let us use the given information about the pole. The object we need to analyze is the pole.



v / The windmills are not closed systems, but angular momentum is being transferred out of them at the same rate it is transferred in, resulting in constant angular momentum. To get an idea of the huge scale of the modern windmill farm, note the sizes of the trucks and trailers.



w / Example 8.

The pole has three forces on it, each of which may also result in a torque: (1) the gravitational force, (2) the cable's force, and (3) the wall's force.

We are free to define an axis of rotation at any point we wish, and it is helpful to define it to lie at the bottom end of the pole, since by that definition the wall's force on the pole is applied at $r = 0$ and thus makes no torque on the pole. This is good, because we don't know what the wall's force on the pole is, and we are not trying to find it.

With this choice of axis, there are two nonzero torques on the pole, a counterclockwise torque from the cable and a clockwise torque from gravity. Choosing to represent counterclockwise torques as positive numbers, and using the equation $|\tau| = r|F| \sin \theta$, we have

$$r_{cable}|F_{cable}| \sin \theta_{cable} - r_{grav}|F_{grav}| \sin \theta_{grav} = 0.$$

A little geometry gives $\theta_{cable} = 90^\circ - \alpha$ and $\theta_{grav} = \alpha$, so

$$r_{cable}|F_{cable}| \sin(90^\circ - \alpha) - r_{grav}|F_{grav}| \sin \alpha = 0.$$

The gravitational force can be considered as acting at the pole's center of mass, i.e., at its geometrical center, so r_{cable} is twice r_{grav} , and we can simplify the equation to read

$$2|F_{cable}| \sin(90^\circ - \alpha) - |F_{grav}| \sin \alpha = 0.$$

These are all quantities we were given, except for α , which is the angle we want to find. To solve for α we need to use the trig identity $\sin(90^\circ - x) = \cos x$,

$$2|F_{cable}| \cos \alpha - |F_{grav}| \sin \alpha = 0,$$

which allows us to find

$$\begin{aligned} \tan \alpha &= 2 \frac{|F_{cable}|}{|F_{grav}|} \\ \alpha &= \tan^{-1} \left(2 \frac{|F_{cable}|}{|F_{grav}|} \right) \\ &= \tan^{-1} \left(2 \times \frac{70 \text{ N}}{98 \text{ N}} \right) \\ &= 55^\circ. \end{aligned}$$

Art!

example 9

▷ The abstract sculpture shown in figure x contains a cube of mass m and sides of length b . The cube rests on top of a cylinder, which is off-center by a distance a . Find the tension in the cable.

▷ There are four forces on the cube: a gravitational force mg , the force F_T from the cable, the upward normal force from the cylinder, F_N , and the horizontal static frictional force from the cylinder, F_s .

The total force on the cube in the vertical direction is zero:

$$F_N - mg = 0.$$

As our axis for defining torques, let's choose the center of the cube. The cable's torque is counterclockwise, the torque due to F_N clockwise. Letting counterclockwise torques be positive, and using the convenient equation $\tau = r_{\perp}F$, we find the equation for the total torque:

$$bF_T - aF_N = 0.$$

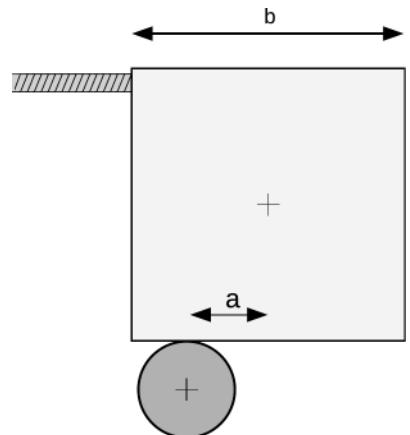
We could also write down the equation saying that the total horizontal force is zero, but that would bring in the cylinder's frictional force on the cube, which we don't know and don't need to find. We already have two equations in the two unknowns F_T and F_N , so there's no need to make it into three equations in three unknowns. Solving the first equation for $F_N = mg$, we then substitute into the second equation to eliminate F_N , and solve for $F_T = (a/b)mg$.

As a check, our result makes sense when $a = 0$; the cube is balanced on the cylinder, so the cable goes slack.

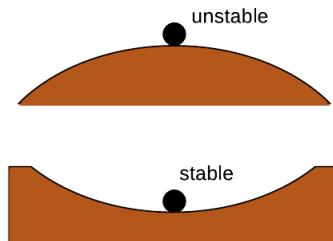
Why is one equilibrium stable and another unstable? Try pushing your own nose to the left or the right. If you push it a millimeter to the left, it responds with a gentle force to the right. If you push it a centimeter to the left, its force on your finger becomes much stronger. The defining characteristic of a stable equilibrium is that the farther the object is moved away from equilibrium, the stronger the force is that tries to bring it back.

The opposite is true for an unstable equilibrium. In the top figure, the ball resting on the round hill theoretically has zero total force on it when it is exactly at the top. But in reality the total force will not be exactly zero, and the ball will begin to move off to one side. Once it has moved, the net force on the ball is greater than it was, and it accelerates more rapidly. In an unstable equilibrium, the farther the object gets from equilibrium, the stronger the force that pushes it farther from equilibrium.

This idea can be rephrased in terms of energy. The difference between the stable and unstable equilibria shown in figure y is that in the stable equilibrium, the energy is at a minimum, and moving



x / Example 9.



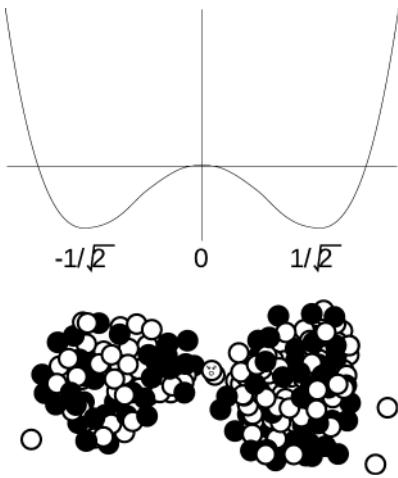
y / Stable and unstable equilibria.



z / The dancer's equilibrium is unstable. If she didn't constantly make tiny adjustments, she would tip over.

to either side of equilibrium will increase it, whereas the unstable equilibrium represents a maximum.

Note that we are using the term “stable” in a weaker sense than in ordinary speech. A domino standing upright is stable in the sense we are using, since it will not spontaneously fall over in response to a sneeze from across the room or the vibration from a passing truck. We would only call it unstable in the technical sense if it could be toppled by *any* force, no matter how small. In everyday usage, of course, it would be considered unstable, since the force required to topple it is so small.



aa / Example 10.

An application of calculus

example 10

- ▷ Nancy Neutron is living in a uranium nucleus that is undergoing fission. Nancy’s nuclear energy as a function of position can be approximated by $U = x^4 - x^2$, where all the units and numerical constants have been suppressed for simplicity. Use calculus to locate the equilibrium points, and determine whether they are stable or unstable.
- ▷ The equilibrium points occur where the U is at a minimum or maximum, and minima and maxima occur where the derivative (which equals minus the force on Nancy) is zero. This derivative is $dU/dx = 4x^3 - 2x$, and setting it equal to zero, we have $x = 0, \pm 1/\sqrt{2}$. Minima occur where the second derivative is positive, and maxima where it is negative. The second derivative is $12x^2 - 2$, which is negative at $x = 0$ (unstable) and positive at $x = \pm 1/\sqrt{2}$ (stable). Interpretation: the graph of U is shaped like a rounded letter ‘W’, with the two troughs representing the two halves of the splitting nucleus. Nancy is going to have to decide which half she wants to go with.

4.1.6 Proof of Kepler’s elliptical orbit law

Kepler determined purely empirically that the planets’ orbits were ellipses, without understanding the underlying reason in terms of physical law. Newton’s proof of this fact based on his laws of motion and law of gravity was considered his crowning achievement both by him and by his contemporaries, because it showed that the same physical laws could be used to analyze both the heavens and the earth. Newton’s proof was very lengthy, but by applying the more recent concepts of conservation of energy and angular momentum we can carry out the proof quite simply and succinctly. This subsection can be skipped without losing the continuity of the text.

The basic idea of the proof is that we want to describe the shape of the planet’s orbit with an equation, and then show that this equation is exactly the one that represents an ellipse. Newton’s original proof had to be very complicated because it was based directly on his laws of motion, which include time as a variable. To make any statement about the shape of the orbit, he had to eliminate time

from his equations, leaving only space variables. But conservation laws tell us that certain things don't change over time, so they have already had time eliminated from them.

There are many ways of representing a curve by an equation, of which the most familiar is $y = ax + b$ for a line in two dimensions. It would be perfectly possible to describe a planet's orbit using an x - y equation like this, but remember that we are applying conservation of angular momentum, and the space variables that occur in the equation for angular momentum are the distance from the axis, r , and the angle between the velocity vector and the r vector, which we will call φ . The planet will have $\varphi = 90^\circ$ when it is moving perpendicular to the r vector, i.e., at the moments when it is at its smallest or greatest distances from the sun. When φ is less than 90° the planet is approaching the sun, and when it is greater than 90° it is receding from it. Describing a curve with an r - φ equation is like telling a driver in a parking lot a certain rule for what direction to steer based on the distance from a certain streetlight in the middle of the lot.

The proof is broken into the three parts for easier digestion. The first part is a simple and intuitively reasonable geometrical fact about ellipses, whose proof we relegate to the caption of figure ac; you will not be missing much if you merely absorb the result without reading the proof.

(1) If we use one of the two foci of an ellipse as an axis for defining the variables r and φ , then the angle between the tangent line and the line drawn to the other focus is the same as φ , i.e., the two angles labeled φ in the figure are in fact equal.

The other two parts form the meat of our proof. We state the results first and then prove them.

(2) A planet, moving under the influence of the sun's gravity with less than the energy required to escape, obeys an equation of the form

$$\sin \varphi = \frac{1}{\sqrt{-pr^2 + qr}},$$

where p and q are positive constants that depend on the planet's energy and angular momentum and p is greater than zero.

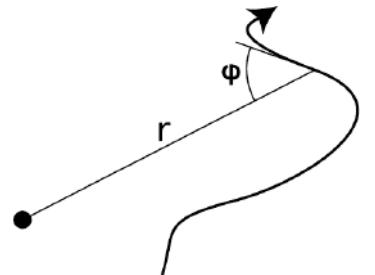
(3) A curve is an ellipse if and only if its r - φ equation is of the form

$$\sin \varphi = \frac{1}{\sqrt{-pr^2 + qr}},$$

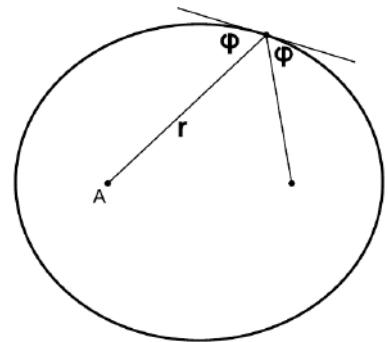
where p and q are positive constants that depend on the size and shape of the ellipse.

Proof of part (2)

The component of the planet's velocity vector that is perpendicular to the \mathbf{r} vector is $v_\perp = v \sin \varphi$, so conservation of angular



ab / Describing a curve by giving φ as a function of r .



ac / Proof that the two angles labeled φ are in fact equal: The definition of an ellipse is that the sum of the distances from the two foci stays constant. If we move a small distance ℓ along the ellipse, then one distance shrinks by an amount $\ell \cos \varphi_1$, while the other grows by $\ell \cos \varphi_2$. These are equal, so $\varphi_1 = \varphi_2$.

momentum tells us that $L = mr v \sin \varphi$ is a constant. Since the planet's mass is a constant, this is the same as the condition

$$rv \sin \varphi = \text{constant}.$$

Conservation of energy gives

$$\frac{1}{2}mv^2 - G\frac{Mm}{r} = \text{constant}.$$

We solve the first equation for v and plug into the second equation to eliminate v . Straightforward algebra then leads to the equation claimed above, with the constant p being positive because of our assumption that the planet's energy is insufficient to escape from the sun, i.e., its total energy is negative.

Proof of part (3)

We define the quantities α , d , and s as shown in figure ad. The law of cosines gives

$$d^2 = r^2 + s^2 - 2rs \cos \alpha.$$

Using $\alpha = 180^\circ - 2\varphi$ and the trigonometric identities $\cos(180^\circ - x) = -\cos x$ and $\cos 2x = 1 - 2 \sin^2 x$, we can rewrite this as

$$d^2 = r^2 + s^2 - 2rs(2 \sin^2 \varphi - 1).$$

Straightforward algebra transforms this into

$$\sin \varphi = \sqrt{\frac{(r+s)^2 - d^2}{4rs}}.$$

Since $r+s$ is constant, the top of the fraction is constant, and the denominator can be rewritten as $4rs = 4r(\text{constant} - r)$, which is equivalent to the desired form.

ad / Quantities referred to in the proof of part (3).

4.2 Rigid-body rotation

4.2.1 Kinematics

When a rigid object rotates, every part of it (every atom) moves in a circle, covering the same angle in the same amount of time, a. Every atom has a different velocity vector, b. Since all the velocities are different, we can't measure the speed of rotation of the top by giving a single velocity. We can, however, specify its speed of rotation consistently in terms of angle per unit time. Let the position of some reference point on the top be denoted by its angle θ , measured in a circle around the axis. For reasons that will become more apparent shortly, we measure all our angles in radians. Then the change in the angular position of any point on the top can be written as $d\theta$, and all parts of the top have the same value of $d\theta$ over a certain time interval dt . We define the angular velocity, ω (Greek omega),

$$\omega = \frac{d\theta}{dt} ,$$

[definition of angular velocity; θ in units of radians]

which is similar to, but not the same as, the quantity ω we defined earlier to describe vibrations. The relationship between ω and t is exactly analogous to that between x and t for the motion of a particle through space.

self-check B

If two different people chose two different reference points on the top in order to define $\theta=0$, how would their $\theta-t$ graphs differ? What effect would this have on the angular velocities? ▷ Answer, p. 1060

The angular velocity has units of radians per second, rad/s. However, radians are not really units at all. The radian measure of an angle is defined, as the length of the circular arc it makes, divided by the radius of the circle. Dividing one length by another gives a unitless quantity, so anything with units of radians is really unitless. We can therefore simplify the units of angular velocity, and call them inverse seconds, s^{-1} .

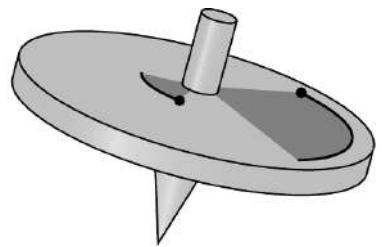
A 78-rpm record

example 11

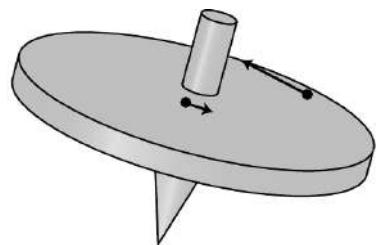
▷ In the early 20th century, the standard format for music recordings was a plastic disk that held a single song and rotated at 78 rpm (revolutions per minute). What was the angular velocity of such a disk?

▷ If we measure angles in units of revolutions and time in units of minutes, then 78 rpm is the angular velocity. Using standard physics units of radians/second, however, we have

$$\frac{78 \text{ revolutions}}{1 \text{ minute}} \times \frac{2\pi \text{ radians}}{1 \text{ revolution}} \times \frac{1 \text{ minute}}{60 \text{ seconds}} = 8.2 \text{ s}^{-1}.$$



a / The two atoms cover the same angle in a given time interval.



b / Their velocity vectors, however, differ in both magnitude and direction.

In the absence of any torque, a rigid body will rotate indefinitely with the same angular velocity. If the angular velocity is changing because of a torque, we define an angular acceleration,

$$\alpha = \frac{d\omega}{dt}, \quad [\text{definition of angular acceleration}]$$

$$x \longleftrightarrow \theta$$

$$v \longleftrightarrow \omega$$

$$a \longleftrightarrow \alpha$$

c / Analogies between rotational and linear quantities.

The symbol is the Greek letter alpha. The units of this quantity are rad/s², or simply s⁻².

The mathematical relationship between ω and θ is the same as the one between v and x , and similarly for α and a . We can thus make a system of analogies, c, and recycle all the familiar kinematic equations for constant-acceleration motion.

The synodic period

example 12

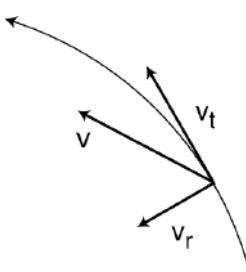
Mars takes nearly twice as long as the Earth to complete an orbit. If the two planets are alongside one another on a certain day, then one year later, Earth will be back at the same place, but Mars will have moved on, and it will take more time for Earth to finish catching up. Angular velocities add and subtract, just as velocity vectors do. If the two planets' angular velocities are ω_1 and ω_2 , then the angular velocity of one relative to the other is $\omega_1 - \omega_2$. The corresponding period, $1/(1/T_1 - 1/T_2)$ is known as the synodic period.

A neutron star

example 13

- ▷ A neutron star is initially observed to be rotating with an angular velocity of 2.0 s^{-1} , determined via the radio pulses it emits. If its angular acceleration is a constant $-1.0 \times 10^{-8} \text{ s}^{-2}$, how many rotations will it complete before it stops? (In reality, the angular acceleration is not always constant; sudden changes often occur, and are referred to as "starquakes!")
- ▷ The equation $v_f^2 - v_i^2 = 2a\Delta x$ can be translated into $\omega_f^2 - \omega_i^2 = 2\alpha\Delta\theta$, giving

$$\begin{aligned}\Delta\theta &= (\omega_f^2 - \omega_i^2)/2\alpha \\ &= 2.0 \times 10^8 \text{ radians} \\ &= 3.2 \times 10^7 \text{ rotations.}\end{aligned}$$



d / We construct a coordinate system that coincides with the location and motion of the moving point of interest at a certain moment.

4.2.2 Relations between angular quantities and motion of a point

It is often necessary to be able to relate the angular quantities to the motion of a particular point on the rotating object. As we develop these, we will encounter the first example where the advantages of radians over degrees become apparent.

The speed at which a point on the object moves depends on both the object's angular velocity ω and the point's distance r from the

axis. We adopt a coordinate system, d , with an inward (radial) axis and a tangential axis. The length of the infinitesimal circular arc ds traveled by the point in a time interval dt is related to $d\theta$ by the definition of radian measure, $d\theta = ds/r$, where positive and negative values of ds represent the two possible directions of motion along the tangential axis. We then have $v_t = ds/dt = r d\theta/dt = \omega r$, or

$$v_t = \omega r. \quad [\text{tangential velocity of a point at a distance } r \text{ from the axis of rotation}]$$

The radial component is zero, since the point is not moving inward or outward,

$$v_r = 0. \quad [\text{radial velocity of a point at a distance } r \text{ from the axis of rotation}]$$

Note that we had to use the definition of radian measure in this derivation. Suppose instead we had used units of degrees for our angles and degrees per second for angular velocities. The relationship between $d\theta_{\text{degrees}}$ and ds is $d\theta_{\text{degrees}} = (360/2\pi)s/r$, where the extra conversion factor of $(360/2\pi)$ comes from that fact that there are 360 degrees in a full circle, which is equivalent to 2π radians. The equation for v_t would then have been $v_t = (2\pi/360)(\omega_{\text{degrees per second}})(r)$, which would have been much messier. Simplicity, then, is the reason for using radians rather than degrees; by using radians we avoid infecting all our equations with annoying conversion factors.

Since the velocity of a point on the object is directly proportional to the angular velocity, you might expect that its acceleration would be directly proportional to the angular acceleration. This is not true, however. Even if the angular acceleration is zero, i.e., if the object is rotating at constant angular velocity, every point on it will have an acceleration vector directed toward the axis, e. As derived on page 213, the magnitude of this acceleration is

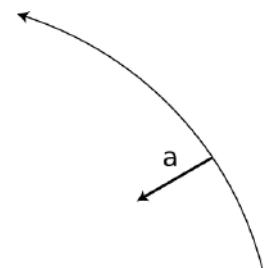
$$a_r = \omega^2 r. \quad [\text{radial acceleration of a point at a distance } r \text{ from the axis}]$$

For the tangential component, any change in the angular velocity $d\omega$ will lead to a change $d\omega \cdot r$ in the tangential velocity, so it is easily shown that

$$a_t = \alpha r. \quad [\text{tangential acceleration of a point at a distance } r \text{ from the axis}]$$

self-check C

Positive and negative signs of ω represent rotation in opposite directions. Why does it therefore make sense physically that ω is raised to the first power in the equation for v_t and to the second power in the one for a_r ? ▷ Answer, p. 1060



e / Even if the rotating object has zero angular acceleration, every point on it has an acceleration towards the center.

Radial acceleration at the surface of the Earth example 14

▷ What is your radial acceleration due to the rotation of the earth if you are at the equator?

▷ At the equator, your distance from the Earth's rotation axis is the same as the radius of the spherical Earth, 6.4×10^6 m. Your angular velocity is

$$\begin{aligned}\omega &= \frac{2\pi \text{ radians}}{1 \text{ day}} \\ &= 7.3 \times 10^{-5} \text{ s}^{-1},\end{aligned}$$

which gives an acceleration of

$$\begin{aligned}a_r &= \omega^2 r \\ &= 0.034 \text{ m/s}^2.\end{aligned}$$

The angular velocity was a very small number, but the radius was a very big number. Squaring a very small number, however, gives a very very small number, so the ω^2 factor “wins,” and the final result is small.

If you're standing on a bathroom scale, this small acceleration is provided by the imbalance between the downward force of gravity and the slightly weaker upward normal force of the scale on your foot. The scale reading is therefore a little lower than it should be.

4.2.3 Dynamics

If we want to connect all this kinematics to anything dynamical, we need to see how it relates to torque and angular momentum. Our strategy will be to tackle angular momentum first, since angular momentum relates to motion, and to use the additive property of angular momentum: the angular momentum of a system of particles equals the sum of the angular momenta of all the individual particles. The angular momentum of one particle within our rigidly rotating object, $L = mv_{\perp}r$, can be rewritten as $L = r p \sin \theta$, where r and p are the magnitudes of the particle's \mathbf{r} and momentum vectors, and θ is the angle between these two vectors. (The \mathbf{r} vector points outward perpendicularly from the axis to the particle's position in space.) In rigid-body rotation the angle θ is 90° , so we have simply $L = rp$. Relating this to angular velocity, we have $L = rp = (r)(mv) = (r)(m\omega r) = mr^2\omega$. The particle's contribution to the total angular momentum is proportional to ω , with a proportionality constant mr^2 . We refer to mr^2 as the particle's contribution to the object's total *moment of inertia*, I , where “moment” is used in the sense of “important,” as in “momentous” — a bigger value of I tells us the particle is more important for determining the total angular momentum. The total moment of inertia

is

$$I = \sum m_i r_i^2, \quad [\text{definition of the moment of inertia;}]$$

for rigid-body rotation in a plane; r is the distance from the axis, measured perpendicular to the axis]

The angular momentum of a rigidly rotating body is then

$$L = I\omega. \quad [\text{angular momentum of rigid-body rotation in a plane}]$$

Since torque is defined as dL/dt , and a rigid body has a constant moment of inertia, we have $\tau = dL/dt = I d\omega/dt = I\alpha$,

$$\tau = I\alpha, \quad [\text{relationship between torque and angular acceleration for rigid-body rotation in a plane}]$$

which is analogous to $F = ma$.

The complete system of analogies between linear motion and rigid-body rotation is given in figure f.

A barbell

example 15

▷ The barbell shown in figure g consists of two small, dense, massive balls at the ends of a very light rod. The balls have masses of 2.0 kg and 1.0 kg, and the length of the rod is 3.0 m. Find the moment of inertia of the rod (1) for rotation about its center of mass, and (2) for rotation about the center of the more massive ball.

▷ (1) The ball's center of mass lies 1/3 of the way from the greater mass to the lesser mass, i.e., 1.0 m from one and 2.0 m from the other. Since the balls are small, we approximate them as if they were two pointlike particles. The moment of inertia is

$$\begin{aligned} I &= (2.0 \text{ kg})(1.0 \text{ m})^2 + (1.0 \text{ kg})(2.0 \text{ m})^2 \\ &= 2.0 \text{ kg}\cdot\text{m}^2 + 4.0 \text{ kg}\cdot\text{m}^2 \\ &= 6.0 \text{ kg}\cdot\text{m}^2 \end{aligned}$$

Perhaps counterintuitively, the less massive ball contributes far more to the moment of inertia.

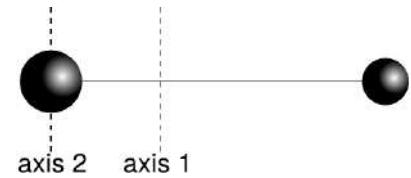
(2) The big ball theoretically contributes a little bit to the moment of inertia, since essentially none of its atoms are exactly at $r=0$. However, since the balls are said to be small and dense, we assume all the big ball's atoms are so close to the axis that we can ignore their small contributions to the total moment of inertia:

$$\begin{aligned} I &= (1.0 \text{ kg})(3.0 \text{ m})^2 \\ &= 9.0 \text{ kg}\cdot\text{m}^2 \end{aligned}$$

This example shows that the moment of inertia depends on the choice of axis. For example, it is easier to wiggle a pen about its center than about one end.

x	\longleftrightarrow	θ
v	\longleftrightarrow	ω
a	\longleftrightarrow	α
m	\longleftrightarrow	I
p	\longleftrightarrow	L
F	\longleftrightarrow	τ

f / Analogies between rotational and linear quantities.



g / Example 15

The parallel axis theorem

example 16

▷ Generalizing the previous example, suppose we pick any axis parallel to axis 1, but offset from it by a distance h . Part (2) of the previous example then corresponds to the special case of $h = -1.0\text{ m}$ (negative being to the left). What is the moment of inertia about this new axis?

▷ The big ball's distance from the new axis is $(1.0\text{ m})+h$, and the small one's is $(2.0\text{ m})-h$. The new moment of inertia is

$$\begin{aligned}I &= (2.0\text{ kg})[(1.0\text{ m})+h]^2 + (1.0\text{ kg})[(2.0\text{ m})-h]^2 \\&= 6.0\text{ kg}\cdot\text{m}^2 + (4.0\text{ kg}\cdot\text{m})h - (4.0\text{ kg}\cdot\text{m})h + (3.0\text{ kg})h^2.\end{aligned}$$

The constant term is the same as the moment of inertia about the center-of-mass axis, the first-order terms cancel out, and the third term is just the total mass multiplied by h^2 . The interested reader will have no difficulty in generalizing this to any set of particles (problem 38, p. 302), resulting in the parallel axis theorem: If an object of total mass M rotates about a line at a distance h from its center of mass, then its moment of inertia equals $I_{cm} + Mh^2$, where I_{cm} is the moment of inertia for rotation about a parallel line through the center of mass.

Scaling of the moment of inertia

example 17

▷ (1) Suppose two objects have the same mass and the same shape, but one is less dense, and larger by a factor k . How do their moments of inertia compare?

(2) What if the densities are equal rather than the masses?

▷ (1) This is like increasing all the distances between atoms by a factor k . All the r 's become greater by this factor, so the moment of inertia is increased by a factor of k^2 .

(2) This introduces an increase in mass by a factor of k^3 , so the moment of inertia of the bigger object is greater by a factor of k^5 .

4.2.4 Iterated integrals

In various places in this book, starting with subsection 4.2.5, we'll come across integrals stuck inside other integrals. These are known as iterated integrals, or double integrals, triple integrals, etc. Similar concepts crop up all the time even when you're not doing calculus, so let's start by imagining such an example. Suppose you want to count how many squares there are on a chess board, and you don't know how to multiply eight times eight. You could start from the upper left, count eight squares across, then continue with the second row, and so on, until you have counted every square, giving the result of 64. In slightly more formal mathematical language, we could write the following recipe: for each row, r , from 1 to 8, consider the columns, c , from 1 to 8, and add one to the count for

each one of them. Using the sigma notation, this becomes

$$\sum_{r=1}^8 \sum_{c=1}^8 1.$$

If you're familiar with computer programming, then you can think of this as a sum that could be calculated using a loop nested inside another loop. To evaluate the result (again, assuming we don't know how to multiply, so we have to use brute force), we can first evaluate the inside sum, which equals 8, giving

$$\sum_{r=1}^8 8.$$

Notice how the “dummy” variable c has disappeared. Finally we do the outside sum, over r , and find the result of 64.

Now imagine doing the same thing with the pixels on a TV screen. The electron beam sweeps across the screen, painting the pixels in each row, one at a time. This is really no different than the example of the chess board, but because the pixels are so small, you normally think of the image on a TV screen as continuous rather than discrete. This is the idea of an integral in calculus. Suppose we want to find the area of a rectangle of width a and height b , and we don't know that we can just multiply to get the area ab . The brute force way to do this is to break up the rectangle into a grid of infinitesimally small squares, each having width dx and height dy , and therefore the infinitesimal area $dA = dx dy$. For convenience, we'll imagine that the rectangle's lower left corner is at the origin. Then the area is given by this integral:

$$\begin{aligned} \text{area} &= \int_{y=0}^b \int_{x=0}^a dA \\ &= \int_{y=0}^b \int_{x=0}^a dx dy \end{aligned}$$

Notice how the leftmost integral sign, over y , and the rightmost differential, dy , act like bookends, or the pieces of bread on a sandwich. Inside them, we have the integral sign that runs over x , and the differential dx that matches it on the right. Finally, on the innermost layer, we'd normally have the thing we're integrating, but here's it's 1, so I've omitted it. Writing the lower limits of the integrals with $x =$ and $y =$ helps to keep it straight which integral goes with which

differential. The result is

$$\begin{aligned}
 \text{area} &= \int_{y=0}^b \int_{x=0}^a dA \\
 &= \int_{y=0}^b \int_{x=0}^a dx dy \\
 &= \int_{y=0}^b \left(\int_{x=0}^a dx \right) dy \\
 &= \int_{y=0}^b a dy \\
 &= a \int_{y=0}^b dy \\
 &= ab.
 \end{aligned}$$

Area of a triangle

example 18

- ▷ Find the area of a 45-45-90 right triangle having legs a .
- ▷ Let the triangle's hypotenuse run from the origin to the point (a, a) , and let its legs run from the origin to $(0, a)$, and then to (a, a) . In other words, the triangle sits on top of its hypotenuse. Then the integral can be set up the same way as the one before, but for a particular value of y , values of x only run from 0 (on the y axis) to y (on the hypotenuse). We then have

$$\begin{aligned}
 \text{area} &= \int_{y=0}^a \int_{x=0}^y dA \\
 &= \int_{y=0}^a \int_{x=0}^y dx dy \\
 &= \int_{y=0}^a \left(\int_{x=0}^y dx \right) dy \\
 &= \int_{y=0}^a y dy \\
 &= \frac{1}{2} a^2
 \end{aligned}$$

Note that in this example, because the upper end of the x values depends on the value of y , it makes a difference which order we do the integrals in. The x integral has to be on the inside, and we have to do it first.

Volume of a cube

example 19

- ▷ Find the volume of a cube with sides of length a .
- ▷ This is a three-dimensional example, so we'll have integrals nested three deep, and the thing we're integrating is the volume $dV = dx dy dz$.

$$\begin{aligned}
 \text{volume} &= \int_{z=0}^a \int_{y=0}^a \int_{x=0}^a dx dy dz \\
 &= \int_{z=0}^a \int_{y=0}^a ady dz \\
 &= a \int_{z=0}^a \int_{y=0}^a dy dz \\
 &= a \int_{z=0}^a adz \\
 &= a^3
 \end{aligned}$$

Area of a circle

example 20

- ▷ Find the area of a circle.
- ▷ To make it easy, let's find the area of a semicircle and then double it. Let the circle's radius be r , and let it be centered on the origin and bounded below by the x axis. Then the curved edge is given by the equation $r^2 = x^2 + y^2$, or $y = \sqrt{r^2 - x^2}$. Since the y integral's limit depends on x , the x integral has to be on the outside. The area is

$$\begin{aligned}
 \text{area} &= \int_{x=-r}^r \int_{y=0}^{\sqrt{r^2-x^2}} dy dx \\
 &= \int_{x=-r}^r \sqrt{r^2 - x^2} dx \\
 &= r \int_{x=-r}^r \sqrt{1 - (x/r)^2} dx.
 \end{aligned}$$

Substituting $u = x/r$,

$$\text{area} = r^2 \int_{u=-1}^1 \sqrt{1 - u^2} du$$

The definite integral equals π , as you can find using a trig substitution or simply by looking it up in a table, and the result is, as expected, $\pi r^2/2$ for the area of the semicircle.

4.2.5 Finding moments of inertia by integration

When calculating the moment of inertia of an ordinary-sized object with perhaps 10^{26} atoms, it would be impossible to do an actual sum over atoms, even with the world's fastest supercomputer. Calculus, however, offers a tool, the integral, for breaking a sum down to infinitely many small parts. If we don't worry about the existence of atoms, then we can use an integral to compute a moment

of inertia as if the object was smooth and continuous throughout, rather than granular at the atomic level. Of course this granularity typically has a negligible effect on the result unless the object is itself an individual molecule. This subsection consists of three examples of how to do such a computation, at three distinct levels of mathematical complication.

Moment of inertia of a thin rod

What is the moment of inertia of a thin rod of mass M and length L about a line perpendicular to the rod and passing through its center? We generalize the discrete sum

$$I = \sum m_i r_i^2$$

to a continuous one,

$$\begin{aligned} I &= \int r^2 dm \\ &= \int_{-L/2}^{L/2} x^2 \frac{M}{L} dx \quad [r = |x|, \text{ so } r^2 = x^2] \\ &= \frac{1}{12} ML^2 \end{aligned}$$

In this example the object was one-dimensional, which made the math simple. The next example shows a strategy that can be used to simplify the math for objects that are three-dimensional, but possess some kind of symmetry.

Moment of inertia of a disk

What is the moment of inertia of a disk of radius b , thickness t , and mass M , for rotation about its central axis?

We break the disk down into concentric circular rings of thickness dr . Since all the mass in a given circular slice has essentially the same value of r (ranging only from r to $r + dr$), the slice's contribution to the total moment of inertia is simply $r^2 dm$. We then have

$$\begin{aligned} I &= \int r^2 dm \\ &= \int r^2 \rho dV, \end{aligned}$$

where $V = \pi b^2 t$ is the total volume, $\rho = M/V = M/\pi b^2 t$ is the density, and the volume of one slice can be calculated as the volume enclosed by its outer surface minus the volume enclosed by its inner surface, $dV = \pi(r + dr)^2 t - \pi r^2 t = 2\pi r t dr$.

$$\begin{aligned} I &= \int_0^b r^2 \frac{M}{\pi b^2 t} 2\pi r t dr \\ &= \frac{1}{2} Mb^2. \end{aligned}$$

In the most general case where there is no symmetry about the rotation axis, we must use iterated integrals, as discussed in subsection 4.2.4. The example of the disk possessed two types of symmetry with respect to the rotation axis: (1) the disk is the same when rotated through any angle about the axis, and (2) all slices perpendicular to the axis are the same. These two symmetries reduced the number of layers of integrals from three to one. The following example possesses only one symmetry, of type (2), and we simply set it up as a triple integral. You may not have seen multiple integrals yet in a math course. If so, just skim this example.

Moment of inertia of a cube

What is the moment of inertia of a cube of side b , for rotation about an axis that passes through its center and is parallel to four of its faces? Let the origin be at the center of the cube, and let x be the rotation axis.

$$\begin{aligned} I &= \int r^2 dm \\ &= \rho \int r^2 dV \\ &= \rho \int_{-b/2}^{b/2} \int_{-b/2}^{b/2} \int_{-b/2}^{b/2} (y^2 + z^2) dx dy dz \\ &= \rho b \int_{-b/2}^{b/2} \int_{-b/2}^{b/2} (y^2 + z^2) dy dz \end{aligned}$$

The fact that the last step is a trivial integral results from the symmetry of the problem. The integrand of the remaining double integral breaks down into two terms, each of which depends on only one of the variables, so we break it into two integrals,

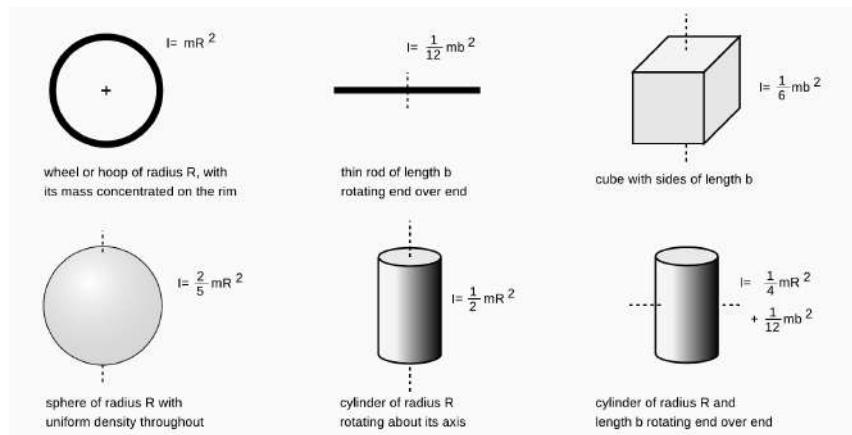
$$I = \rho b \int_{-b/2}^{b/2} \int_{-b/2}^{b/2} y^2 dy dz + \rho b \int_{-b/2}^{b/2} \int_{-b/2}^{b/2} z^2 dy dz$$

which we know have identical results. We therefore only need to evaluate one of them and double the result:

$$\begin{aligned} I &= 2\rho b \int_{-b/2}^{b/2} \int_{-b/2}^{b/2} z^2 dy dz \\ &= 2\rho b^2 \int_{-b/2}^{b/2} z^2 dz \\ &= \frac{1}{6} \rho b^5 \\ &= \frac{1}{6} Mb^2 \end{aligned}$$

Figure h shows the moments of inertia of some shapes, which were evaluated with techniques like these.

h / Moments of inertia of some geometric shapes



The hammer throw

example 21

▷ In the men's Olympic hammer throw, a steel ball of radius 6.1 cm is swung on the end of a wire of length 1.22 m. What fraction of the ball's angular momentum comes from its rotation, as opposed to its motion through space?

▷ It's always important to solve problems symbolically first, and plug in numbers only at the end, so let the radius of the ball be b , and the length of the wire ℓ . If the time the ball takes to go once around the circle is T , then this is also the time it takes to revolve once around its own axis. Its speed is $v = 2\pi\ell/T$, so its angular momentum due to its motion through space is $mv\ell = 2\pi m\ell^2/T$. Its angular momentum due to its rotation around its own center is $(4\pi/5)mb^2/T$. The ratio of these two angular momenta is $(2/5)(b/\ell)^2 = 1.0 \times 10^{-3}$. The angular momentum due to the ball's spin is extremely small.

Toppling a rod

example 22

▷ A rod of length b and mass m stands upright. We want to strike the rod at the bottom, causing it to fall and land flat. Find the momentum, p , that should be delivered, in terms of m , b , and g . Can this really be done without having the rod scrape on the floor?

▷ This is a nice example of a question that can very nearly be answered based only on units. Since the three variables, m , b , and g , all have different units, they can't be added or subtracted. The only way to combine them mathematically is by multiplication or division. Multiplying one of them by itself is exponentiation, so in general we expect that the answer must be of the form

$$p = Am^j b^k g^l,$$

where A , j , k , and l are unitless constants. The result has to have units of $\text{kg}\cdot\text{m/s}$. To get kilograms to the first power, we need

$$j = 1,$$

i / Example 22.



meters to the first power requires

$$k + l = 1,$$

and seconds to the power -1 implies

$$l = 1/2.$$

We find $j = 1$, $k = 1/2$, and $l = 1/2$, so the solution must be of the form

$$p = Am\sqrt{bg}.$$

Note that no physics was required!

Consideration of units, however, won't help us to find the unitless constant A . Let t be the time the rod takes to fall, so that $(1/2)gt^2 = b/2$. If the rod is going to land exactly on its side, then the number of revolutions it completes while in the air must be $1/4$, or $3/4$, or $5/4$, \dots , but all the possibilities greater than $1/4$ would cause the head of the rod to collide with the floor prematurely. The rod must therefore rotate at a rate that would cause it to complete a full rotation in a time $T = 4t$, and it has angular momentum $L = (\pi/6)mb^2/T$.

The momentum lost by the object striking the rod is p , and by conservation of momentum, this is the amount of momentum, in the horizontal direction, that the rod acquires. In other words, the rod will fly forward a little. However, this has no effect on the solution to the problem. More importantly, the object striking the rod loses angular momentum $bp/2$, which is also transferred to the rod. Equating this to the expression above for L , we find $p = (\pi/12)m\sqrt{bg}$.

Finally, we need to know whether this can really be done without having the foot of the rod scrape on the floor. The figure shows that the answer is no for this rod of finite width, but it appears that the answer would be yes for a sufficiently thin rod. This is analyzed further in homework problem 37 on page 301.

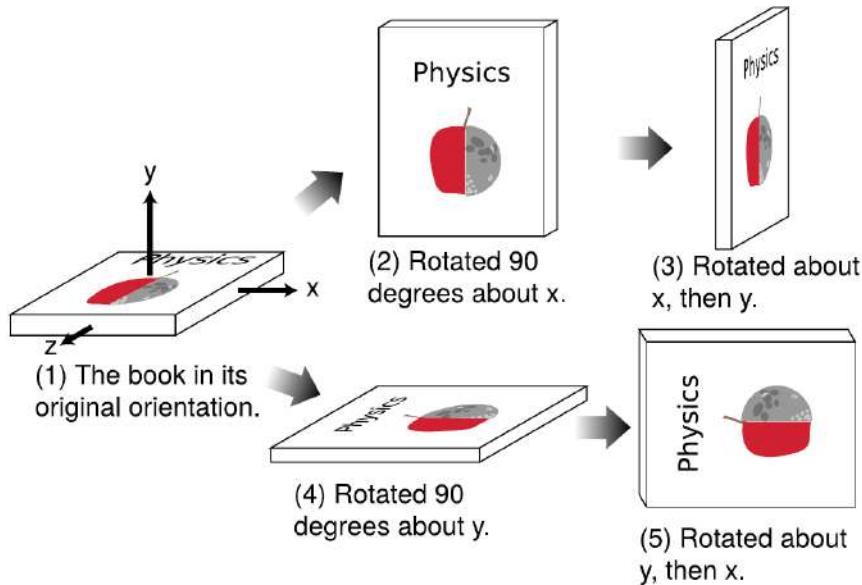
4.3 Angular momentum in three dimensions

Conservation of angular momentum produces some surprising phenomena when extended to three dimensions. Try the following experiment, for example. Take off your shoe, and toss it in to the air, making it spin along its long (toe-to-heel) axis. You should observe a nice steady pattern of rotation. The same happens when you spin the shoe about its shortest (top-to-bottom) axis. But something unexpected happens when you spin it about its third (left-to-right) axis, which is intermediate in length between the other two. Instead of a steady pattern of rotation, you will observe something more complicated, with the shoe changing its orientation with respect to the rotation axis.

4.3.1 Rigid-body kinematics in three dimensions

How do we generalize rigid-body kinematics to three dimensions? When we wanted to generalize the kinematics of a moving particle to three dimensions, we made the numbers r , v , and a into vectors \mathbf{r} , \mathbf{v} , and \mathbf{a} . This worked because these quantities all obeyed the same laws of vector addition. For instance, one of the laws of vector addition is that, just like addition of numbers, vector addition gives the same result regardless of the order of the two quantities being added. Thus you can step sideways 1 m to the right and then step forward 1 m, and the end result is the same as if you stepped forward first and then to the side. In other words, it didn't matter whether you took $\Delta\mathbf{r}_1 + \Delta\mathbf{r}_2$ or $\Delta\mathbf{r}_2 + \Delta\mathbf{r}_1$. In math this is called the commutative property of addition.

a / Performing the rotations in one order gives one result, 3, while reversing the order gives a different result, 5.



Angular motion, unfortunately doesn't have this property, as shown in figure a. Doing a rotation about the x axis and then

about y gives one result, while doing them in the opposite order gives a different result. These operations don't "commute," i.e., it makes a difference what order you do them in.

This means that there is in general no possible way to construct a $\Delta\theta$ vector. However, if you try doing the operations shown in figure a using small rotation, say about 10 degrees instead of 90, you'll find that the result is nearly the same regardless of what order you use; small rotations are very nearly commutative. Not only that, but the result of the two 10-degree rotations is about the same as a single, somewhat larger, rotation about an axis that lies symmetrically between the x and y axes at 45 degree angles to each one. This is exactly what we would expect if the two small rotations did act like vectors whose directions were along the axis of rotation. We therefore define a $d\theta$ vector whose magnitude is the amount of rotation in units of radians, and whose direction is along the axis of rotation. Actually this definition is ambiguous, because there it could point in either direction along the axis. We therefore use a right-hand rule as shown in figure b to define the direction of the $d\theta$ vector, and the ω vector, $\omega = d\theta/dt$, based on it. Aliens on planet Tammyfaye may decide to define it using their left hands rather than their right, but as long as they keep their scientific literature separate from ours, there is no problem. When entering a physics exam, always be sure to write a large warning note on your left hand in magic marker so that you won't be tempted to use it for the right-hand rule while keeping your pen in your right.

self-check D

Use the right-hand rule to determine the directions of the ω vectors in each rotation shown in figures a/1 through a/5. \triangleright Answer, p. 1060

Because the vector relationships among $d\theta$, ω , and α are strictly analogous to the ones involving dr , v , and a (with the proviso that we avoid describing large rotations using $\Delta\theta$ vectors), any operation in $r\text{-}v\text{-}a$ vector kinematics has an exact analog in $\theta\text{-}\omega\text{-}\alpha$ kinematics.

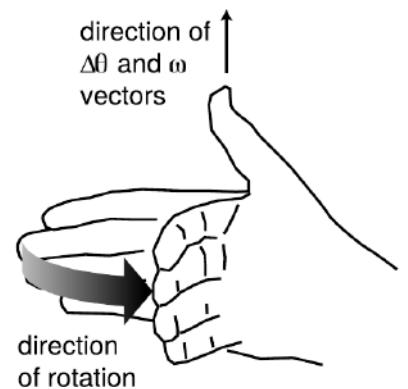
Result of successive 10-degree rotations

example 23

- \triangleright What is the result of two successive (positive) 10-degree rotations about the x and y axes? That is, what single rotation about a single axis would be equivalent to executing these in succession?
- \triangleright The result is only going to be approximate, since 10 degrees is not an infinitesimally small angle, and we are not told in what order the rotations occur. To some approximation, however, we can add the $\Delta\theta$ vectors in exactly the same way we would add Δr vectors, so we have

$$\begin{aligned}\Delta\theta &\approx \Delta\theta_1 + \Delta\theta_2 \\ &\approx (10 \text{ degrees})\hat{x} + (10 \text{ degrees})\hat{y}.\end{aligned}$$

This is a vector with a magnitude of $\sqrt{(10 \text{ deg})^2 + (10 \text{ deg})^2} =$



b / The right-hand rule for associating a vector with a direction of rotation.

14 deg, and it points along an axis midway between the x and y axes.

4.3.2 Angular momentum in three dimensions

The vector cross product

In order to expand our system of three-dimensional kinematics to include dynamics, we will have to generalize equations like $v_t = \omega r$, $\tau = rF \sin \theta_{rF}$, and $L = rp \sin \theta_{rp}$, each of which involves three quantities that we have either already defined as vectors or that we want to redefine as vectors. Although the first one appears to differ from the others in its form, it could just as well be rewritten as $v_t = \omega r \sin \theta_{\omega r}$, since $\theta_{\omega r} = 90^\circ$, and $\sin \theta_{\omega r} = 1$.

It thus appears that we have discovered something general about the physically useful way to relate three vectors in a multiplicative way: the magnitude of the result always seems to be proportional to the product of the magnitudes of the two vectors being “multiplied,” and also to the sine of the angle between them.

Is this pattern just an accident? Actually the sine factor has a very important physical property: it goes to zero when the two vectors are parallel. This is a Good Thing. The generalization of angular momentum into a three-dimensional vector, for example, is presumably going to describe not just the clockwise or counterclockwise nature of the motion but also from which direction we would have to view the motion so that it was clockwise or counterclockwise. (A clock’s hands go counterclockwise as seen from behind the clock, and don’t rotate at all as seen from above or to the side.) Now suppose a particle is moving directly away from the origin, so that its \mathbf{r} and \mathbf{p} vectors are parallel. It is not going around the origin from any point of view, so its angular momentum vector had better be zero.

Thinking in a slightly more abstract way, we would expect the angular momentum vector to point perpendicular to the plane of motion, just as the angular velocity vector points perpendicular to the plane of motion. The plane of motion is the plane containing both \mathbf{r} and \mathbf{p} , if we place the two vectors tail-to-tail. But if \mathbf{r} and \mathbf{p} are parallel and are placed tail-to-tail, then there are infinitely many planes containing them both. To pick one of these planes in preference to the others would violate the symmetry of space, since they should all be equally good. Thus the zero-if-parallel property is a necessary consequence of the underlying symmetry of the laws of physics.

The following definition of a kind of vector multiplication is consistent with everything we’ve seen so far, and on p. 1027 we’ll prove that the definition is unique, i.e., if we believe in the symmetry of space, it is essentially the only way of defining the multiplication of

two vectors to produce a third vector:

Definition of the vector cross product:

The cross product $\mathbf{A} \times \mathbf{B}$ of two vectors \mathbf{A} and \mathbf{B} is defined as follows:

- (1) Its magnitude is defined by $|\mathbf{A} \times \mathbf{B}| = |\mathbf{A}||\mathbf{B}| \sin \theta_{AB}$, where θ_{AB} is the angle between \mathbf{A} and \mathbf{B} when they are placed tail-to-tail.
- (2) Its direction is along the line perpendicular to both \mathbf{A} and \mathbf{B} . Of the two such directions, it is the one that obeys the right-hand rule shown in figure c.

The name “cross product” refers to the symbol, and distinguishes it from the dot product, which acts on two vectors but produces a scalar.

Although the vector cross-product has nearly all the properties of numerical multiplication, e.g., $\mathbf{A} \times (\mathbf{B} + \mathbf{C}) = \mathbf{A} \times \mathbf{B} + \mathbf{A} \times \mathbf{C}$, it lacks the usual property of commutativity. Try applying the right-hand rule to find the direction of the vector cross product $\mathbf{B} \times \mathbf{A}$ using the two vectors shown in the figure. This requires starting with a flattened hand with the four fingers pointing along \mathbf{B} , and then curling the hand so that the fingers point along \mathbf{A} . The only possible way to do this is to point your thumb toward the floor, in the opposite direction. Thus for the vector cross product we have

$$\mathbf{A} \times \mathbf{B} = -\mathbf{B} \times \mathbf{A},$$

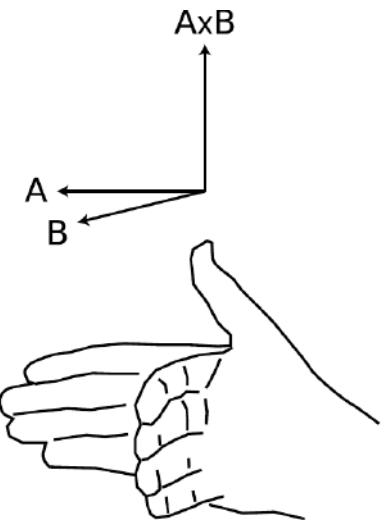
a property known as anticommutativity. The vector cross product is also not associative, i.e., $\mathbf{A} \times (\mathbf{B} \times \mathbf{C})$ is usually not the same as $(\mathbf{A} \times \mathbf{B}) \times \mathbf{C}$.

A geometric interpretation of the cross product, d, is that if both \mathbf{A} and \mathbf{B} are vectors with units of distance, then the magnitude of their cross product can be interpreted as the area of the parallelogram they form when placed tail-to-tail.

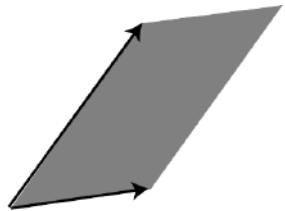
A useful expression for the components of the vector cross product in terms of the components of the two vectors being multiplied is as follows:

$$\begin{aligned}(\mathbf{A} \times \mathbf{B})_x &= A_y B_z - B_y A_z \\(\mathbf{A} \times \mathbf{B})_y &= A_z B_x - B_z A_x \\(\mathbf{A} \times \mathbf{B})_z &= A_x B_y - B_x A_y\end{aligned}$$

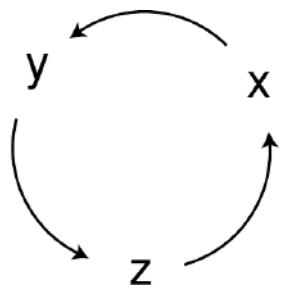
I'll prove later that these expressions are equivalent to the previous definition of the cross product. Although they may appear formidable, they have a simple structure: the subscripts on the right are the other two besides the one on the left, and each equation is related to the preceding one by a cyclic change in the subscripts, e. If the subscripts were not treated in some completely symmetric



c / The right-hand rule for the direction of the vector cross product.



d / The magnitude of the cross product is the area of the shaded parallelogram.



e / A cyclic change in the x , y , and z subscripts.

manner like this, then the definition would provide some way to distinguish one axis from another, which would violate the symmetry of space.

self-check E

Show that the component equations are consistent with the rule $\mathbf{A} \times \mathbf{B} = -\mathbf{B} \times \mathbf{A}$.
 ▷ Answer, p. 1060

Angular momentum in three dimensions

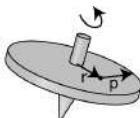
In terms of the vector cross product, we have:

$$\mathbf{v} = \boldsymbol{\omega} \times \mathbf{r}$$

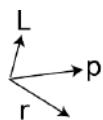
$$\mathbf{L} = \mathbf{r} \times \mathbf{p}$$

$$\boldsymbol{\tau} = \mathbf{r} \times \mathbf{F}$$

But wait, how do we know these equations are even correct? For instance, how do we know that the quantity defined by $\mathbf{r} \times \mathbf{p}$ is in fact conserved? Well, just as we saw on page 216 that the dot product is unique (i.e., can only be defined in one way while observing rotational invariance), the cross product is also unique, as proved on page 1027. If $\mathbf{r} \times \mathbf{p}$ was not conserved, then there could not be any generally conserved quantity that would reduce to our old definition of angular momentum in the special case of plane rotation. This doesn't prove conservation of angular momentum — only experiments can prove that — but it does prove that if angular momentum is conserved in three dimensions, there is only one possible way to generalize from two dimensions to three.



f / The position and momentum vectors of an atom in the spinning top.



g / The right-hand rule for the atom's contribution to the angular momentum.

Angular momentum of a spinning top

example 24

As an illustration, we consider the angular momentum of a spinning top. Figures f and g show the use of the vector cross product to determine the contribution of a representative atom to the total angular momentum. Since every other atom's angular momentum vector will be in the same direction, this will also be the direction of the total angular momentum of the top. This happens to be rigid-body rotation, and perhaps not surprisingly, the angular momentum vector is along the same direction as the angular velocity vector.

Three important points are illustrated by this example: (1) When we do the full three-dimensional treatment of angular momentum, the “axis” from which we measure the position vectors is just an arbitrarily chosen point. If this had not been rigid-body rotation, we would not even have been able to identify a single line about which every atom circled. (2) Starting from figure f, we had to rearrange the vectors to get them tail-to-tail before applying the right-hand rule. If we had attempted to apply the right-hand rule to figure f, the direction of the result would have been exactly the opposite of the correct answer. (3) The equation $\mathbf{L} = \mathbf{r} \times \mathbf{p}$ cannot be applied all at once to an entire system of particles. The total

momentum of the top is zero, which would give an erroneous result of zero angular momentum (never mind the fact that \mathbf{r} is not well defined for the top as a whole).

Doing the right-hand rule like this requires some practice. I urge you to make models like g out of rolled up pieces of paper and to practice with the model in various orientations until it becomes natural.

Precession

example 25

Figure h shows a counterintuitive example of the concepts we've been discussing. One expects the torque due to gravity to cause the top to flop down. Instead, the top remains spinning in the horizontal plane, but its axis of rotation starts moving in the direction shown by the shaded arrow. This phenomenon is called precession. Figure i shows that the torque due to gravity is out of the page. (Actually we should add up all the torques on all the atoms in the top, but the qualitative result is the same.) Since torque is the rate of change of angular momentum, $\tau = d\mathbf{L}/dt$, the $\Delta\mathbf{L}$ vector must be in the same direction as the torque (division by a positive scalar doesn't change the direction of the vector). As shown in j, this causes the angular momentum vector to twist in space without changing its magnitude.

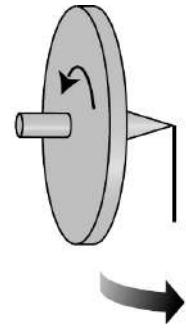
For similar reasons, the Earth's axis precesses once every 26,000 years (although not through a great circle, since the angle between the axis and the force isn't 90 degrees as in figure h). This precession is due to a torque exerted by the moon. If the Earth was a perfect sphere, there could be no precession effect due to symmetry. However, the Earth's own rotation causes it to be slightly flattened (oblate) relative to a perfect sphere, giving it "love handles" on which the moon's gravity can act. The moon's gravity on the nearer side of the equatorial bulge is stronger, so the torques do not cancel out perfectly. Presently the earth's axis very nearly lines up with the star Polaris, but in 12,000 years, the pole star will be Vega instead.

The frisbee

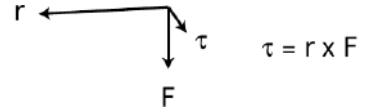
example 26

The flow of the air over a flying frisbee generates lift, and the lift at the front and back of the frisbee isn't necessarily balanced. If you throw a frisbee without rotating it, as if you were shooting a basketball with two hands, you'll find that it pitches, i.e., its nose goes either up or down. When I do this with my frisbee, it goes nose down, which apparently means that the lift at the back of the disc is greater than the lift at the front. The two torques are unbalanced, resulting in a total torque that points to the left.

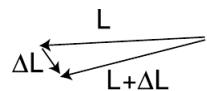
The way you actually throw a frisbee is with one hand, putting a lot of spin on it. If you throw backhand, which is how most people first learn to do it, the angular momentum vector points down (assuming you're right-handed). On my frisbee, the aerodynamic



h / A top is supported at its tip by a pinhead. (More practical devices to demonstrate this would use a double bearing.)



i / The torque made by gravity is in the horizontal plane.



j / The $\Delta\mathbf{L}$ vector is in the same direction as the torque, out of the page.

torque to the left would therefore tend to make the angular momentum vector precess in the clockwise direction as seen by the thrower. This would cause the disc to roll to the right, and therefore follow a curved trajectory. Some specialized discs, used in the sport of disc golf, are actually designed intentionally to show this behavior; they're known as "understable" discs. However, the typical frisbee that most people play with is designed to be stable: as the disc rolls to one side, the airflow around it is altered in way that tends to bring the disc back into level flight. Such a disc will therefore tend to fly in a straight line, provided that it is thrown with enough angular momentum.

r	4	5	0
F	1	2	3

k / Example 27.

Finding a cross product by components example 27

- ▷ What is the torque produced by a force given by $\hat{x} + 2\hat{y} + 3\hat{z}$ (in units of Newtons) acting on a point whose radius vector is $4\hat{x} + 5\hat{y}$ (in meters)?

- ▷ It's helpful to make a table of the components as shown in the figure. The results are

$$\begin{aligned}\tau_x &= r_y F_z - F_y r_z = 15 \text{ N}\cdot\text{m} \\ \tau_y &= r_z F_x - F_z r_x = -12 \text{ N}\cdot\text{m} \\ \tau_z &= r_x F_y - F_x r_y = 3 \text{ N}\cdot\text{m}\end{aligned}$$

Torque and angular momentum example 28

In this example, we prove explicitly the consistency of the equations involving torque and angular momentum that we proved above based purely on symmetry. Starting from the definition of torque, we have

$$\begin{aligned}\tau &= \frac{d\mathbf{L}}{dt} \\ &= \frac{d}{dt} \sum_i \mathbf{r}_i \times \mathbf{p}_i \\ &= \sum_i \frac{d}{dt} (\mathbf{r}_i \times \mathbf{p}_i).\end{aligned}$$

The derivative of a cross product can be evaluated in the same way as the derivative of an ordinary scalar product:

$$\tau = \sum_i \left[\left(\frac{d\mathbf{r}_i}{dt} \times \mathbf{p}_i \right) + \left(\mathbf{r}_i \times \frac{d\mathbf{p}_i}{dt} \right) \right]$$

The first term is zero for each particle, since the velocity vector is parallel to the momentum vector. The derivative appearing in the second term is the force acting on the particle, so

$$\tau = \sum_i \mathbf{r}_i \times \mathbf{F}_i,$$

which is the relationship we set out to prove.

4.3.3 Rigid-body dynamics in three dimensions

The student who is not madly in love with mathematics may wish to skip the rest of this section after absorbing the statement that, for a typical, asymmetric object, the angular momentum vector and the angular velocity vector need not be parallel. That is, only for a body that possesses symmetry about the rotation axis is it true that $\mathbf{L} = I\boldsymbol{\omega}$ (the rotational equivalent of $\mathbf{p} = m\mathbf{v}$) for some scalar I .

Let's evaluate the angular momentum of a rigidly rotating system of particles:

$$\begin{aligned}\mathbf{L} &= \sum_i \mathbf{r}_i \times \mathbf{p}_i \\ &= \sum_i m_i \mathbf{r}_i \times \mathbf{v}_i \\ &= \sum_i m_i \mathbf{r}_i \times (\boldsymbol{\omega} \times \mathbf{r}_i)\end{aligned}$$

An important mathematical skill is to know when to give up and back off. This is a complicated expression, and there is no reason to expect it to simplify and, for example, take the form of a scalar multiplied by $\boldsymbol{\omega}$. Instead we examine its general characteristics. If we expanded it using the equation that gives the components of a vector cross product, every term would have one of the $\boldsymbol{\omega}$ components raised to the first power, multiplied by a bunch of other stuff. The most general possible form for the result is

$$\begin{aligned}L_x &= I_{xx}\omega_x + I_{xy}\omega_y + I_{xz}\omega_z \\ L_y &= I_{yx}\omega_x + I_{yy}\omega_y + I_{yz}\omega_z \\ L_z &= I_{zx}\omega_x + I_{zy}\omega_y + I_{zz}\omega_z,\end{aligned}$$

which you may recognize as a case of matrix multiplication. The moment of inertia is not a scalar, and not a three-component vector. It is a matrix specified by nine numbers, called its matrix elements.

The elements of the moment of inertia matrix will depend on our choice of a coordinate system. In general, there will be some special coordinate system, in which the matrix has a simple diagonal form:

$$\begin{aligned}L_x &= I_{xx}\omega_x \\ L_y &= I_{yy}\omega_y \\ L_z &= I_{zz}\omega_z.\end{aligned}$$

The three special axes that cause this simplification are called the principal axes of the object, and the corresponding coordinate system is the principal axis system. For symmetric shapes such as a rectangular box or an ellipsoid, the principal axes lie along the intersections of the three symmetry planes, but even an asymmetric body has principal axes.

We can also generalize the plane-rotation equation $K = (1/2)I\omega^2$ to three dimensions as follows:

$$\begin{aligned} K &= \sum_i \frac{1}{2} m_i v_i^2 \\ &= \frac{1}{2} \sum_i m_i (\boldsymbol{\omega} \times \mathbf{r}_i) \cdot (\boldsymbol{\omega} \times \mathbf{r}_i) \end{aligned}$$

We want an equation involving the moment of inertia, and this has some evident similarities to the sum we originally wrote down for the moment of inertia. To massage it into the right shape, we need the vector identity $(\mathbf{A} \times \mathbf{B}) \cdot \mathbf{C} = (\mathbf{B} \times \mathbf{C}) \cdot \mathbf{A}$, which we state without proof. We then write

$$\begin{aligned} K &= \frac{1}{2} \sum_i m_i [\mathbf{r}_i \times (\boldsymbol{\omega} \times \mathbf{r}_i)] \cdot \boldsymbol{\omega} \\ &= \frac{1}{2} \boldsymbol{\omega} \cdot \sum_i m_i \mathbf{r}_i \times (\boldsymbol{\omega} \times \mathbf{r}_i) \\ &= \frac{1}{2} \mathbf{L} \cdot \boldsymbol{\omega} \end{aligned}$$

As a reward for all this hard work, let's analyze the problem of the spinning shoe that I posed at the beginning of the chapter. The three rotation axes referred to there are approximately the principal axes of the shoe. While the shoe is in the air, no external torques are acting on it, so its angular momentum vector must be constant in magnitude and direction. Its kinetic energy is also constant. That's in the room's frame of reference, however. The principal axis frame is attached to the shoe, and tumbles madly along with it. In the principal axis frame, the kinetic energy and the magnitude of the angular momentum stay constant, but the actual direction of the angular momentum need not stay fixed (as you saw in the case of rotation that was initially about the intermediate-length axis). Constant $|\mathbf{L}|$ gives

$$L_x^2 + L_y^2 + L_z^2 = \text{constant.}$$

In the principal axis frame, it's easy to solve for the components of $\boldsymbol{\omega}$ in terms of the components of \mathbf{L} , so we eliminate $\boldsymbol{\omega}$ from the expression $2K = \mathbf{L} \cdot \boldsymbol{\omega}$, giving

$$\frac{1}{I_{xx}} L_x^2 + \frac{1}{I_{yy}} L_y^2 + \frac{1}{I_{zz}} L_z^2 = \text{constant } \#2.$$

The first equation is the equation of a sphere in the three dimensional space occupied by the angular momentum vector, while the second one is the equation of an ellipsoid. The top figure corresponds to the case of rotation about the shortest axis, which has the greatest moment of inertia element. The intersection of the two

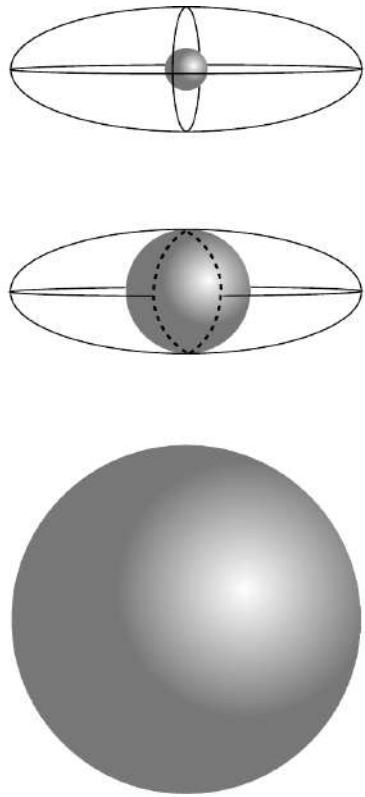
surfaces consists only of the two points at the front and back of the sphere. The angular momentum is confined to one of these points, and can't change its direction, i.e., its orientation with respect to the principal axis system, which is another way of saying that the shoe can't change its orientation with respect to the angular momentum vector. In the bottom figure, the shoe is rotating about the longest axis. Now the angular momentum vector is trapped at one of the two points on the right or left. In the case of rotation about the axis with the intermediate moment of inertia element, however, the intersection of the sphere and the ellipsoid is not just a pair of isolated points but the curve shown with the dashed line. The relative orientation of the shoe and the angular momentum vector can and will change.

One application of the moment of inertia tensor is to video games that simulate car racing or flying airplanes.

One more exotic example has to do with nuclear physics. Although you have probably visualized atomic nuclei as nothing more than featureless points, or perhaps tiny spheres, they are often ellipsoids with one long axis and two shorter, equal ones. Although a spinning nucleus normally gets rid of its angular momentum via gamma ray emission within a period of time on the order of picoseconds, it may happen that a deformed nucleus gets into a state in which has a large angular momentum is along its long axis, which is a very stable mode of rotation. Such states can live for seconds or even years! (There is more to the story — this is the topic on which I wrote my Ph.D. thesis — but the basic insight applies even though the full treatment requires fancy quantum mechanics.)

Our analysis has so far assumed that the kinetic energy of rotation energy can't be converted into other forms of energy such as heat, sound, or vibration. When this assumption fails, then rotation about the axis of least moment of inertia becomes unstable, and will eventually convert itself into rotation about the axis whose moment of inertia is greatest. This happened to the U.S.'s first artificial satellite, Explorer I, launched in 1958. Note the long, floppy antennas, which tended to dissipate kinetic energy into vibration. It had been designed to spin about its minimum-moment-of-inertia axis, but almost immediately, as soon as it was in space, it began spinning end over end. It was nevertheless able to carry out its science mission, which didn't depend on being able to maintain a stable orientation, and it discovered the Van Allen radiation belts.

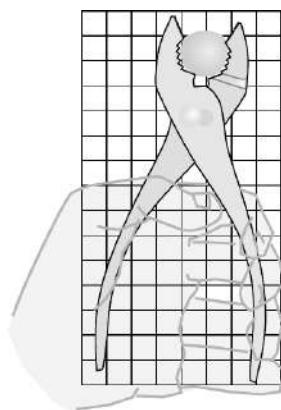
This chapter is summarized on page 1079. Notation and terminology are tabulated on pages 1070-1071.



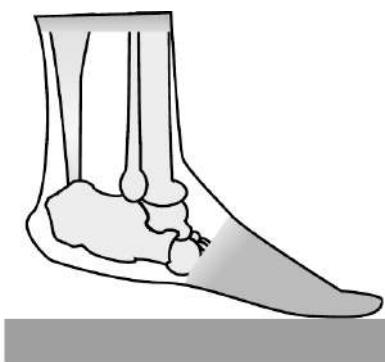
I / Visualizing surfaces of constant energy and angular momentum in L_x - L_y - L_z space.



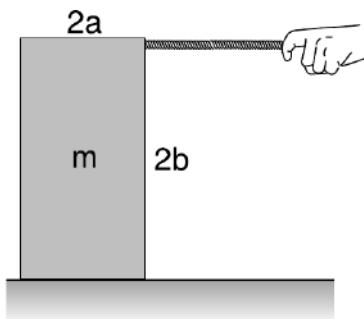
m / The Explorer I satellite.



Problem 1.



Problem 6.



Problem 8.

Problems

The symbols \checkmark , \blacksquare , etc. are explained on page 303.

- 1** The figure shows scale drawing of a pair of pliers being used to crack a nut, with an appropriately reduced centimeter grid. Warning: do not attempt this at home; it is bad manners. If the force required to crack the nut is 300 N, estimate the force required of the person's hand.

\triangleright Solution, p. 1043 \blacksquare

- 2** You are trying to loosen a stuck bolt on your RV using a big wrench that is 50 cm long. If you hang from the wrench, and your mass is 55 kg, what is the maximum torque you can exert on the bolt? \checkmark \blacksquare

- 3** A physical therapist wants her patient to rehabilitate his injured elbow by laying his arm flat on a table, and then lifting a 2.1 kg mass by bending his elbow. In this situation, the weight is 33 cm from his elbow. He calls her back, complaining that it hurts him to grasp the weight. He asks if he can strap a bigger weight onto his arm, only 17 cm from his elbow. How much mass should she tell him to use so that he will be exerting the same torque? (He is raising his forearm itself, as well as the weight.) \checkmark \blacksquare

- 4** An object thrown straight up in the air is momentarily at rest when it reaches the top of its motion. Does that mean that it is in equilibrium at that point? Explain. \blacksquare

- 5** An object is observed to have constant angular momentum. Can you conclude that no torques are acting on it? Explain. [Based on a problem by Serway and Faughn.] \blacksquare

- 6** A person of mass m stands on the ball of one foot. Find the tension in the calf muscle and the force exerted by the shinbones on the bones of the foot, in terms of m , g , a , and b . For simplicity, assume that all the forces are at 90-degree angles to the foot, i.e., neglect the angle between the foot and the floor. \checkmark \blacksquare

- 7** Two pointlike particles have the same momentum vector. Can you conclude that their angular momenta are the same? Explain. [Based on a problem by Serway and Faughn.] \blacksquare

- 8** The box shown in the figure is being accelerated by pulling on it with the rope.

- (a) Assume the floor is frictionless. What is the maximum force that can be applied without causing the box to tip over?

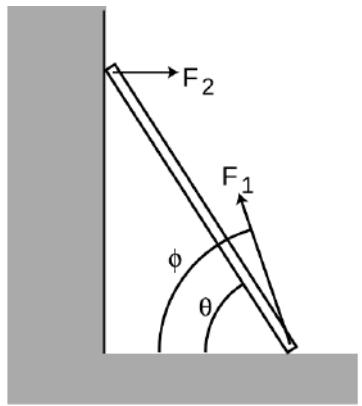
\triangleright Hint, p. 1035 \checkmark

- (b) Repeat part a, but now let the coefficient of friction be μ . \checkmark

- (c) What happens to your answer to part b when the box is sufficiently tall? How do you interpret this? \blacksquare

9 A uniform ladder of mass m and length ℓ leans against a smooth wall, making an angle θ with respect to the ground. The dirt exerts a normal force and a frictional force on the ladder, producing a force vector with magnitude F_1 at an angle ϕ with respect to the ground. Since the wall is smooth, it exerts only a normal force on the ladder; let its magnitude be F_2 .

- (a) Explain why ϕ must be greater than θ . No math is needed.
- (b) Choose any numerical values you like for m and ℓ , and show that the ladder can be in equilibrium (zero torque and zero total force vector) for $\theta=45.00^\circ$ and $\phi=63.43^\circ$. ■



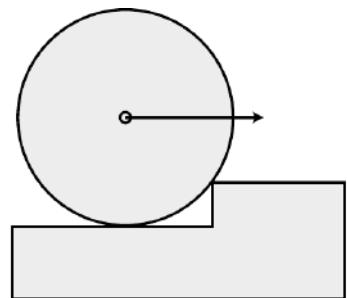
Problems 9 and 10.

10 Continuing problem 9, find an equation for ϕ in terms of θ , and show that m and L do not enter into the equation. Do not assume any numerical values for any of the variables. You will need the trig identity $\sin(a - b) = \sin a \cos b - \sin b \cos a$. (As a numerical check on your result, you may wish to check that the angles given in problem 9b satisfy your equation.) ✓ ■

11 (a) Find the minimum horizontal force which, applied at the axle, will pull a wheel over a step. Invent algebra symbols for whatever quantities you find to be relevant, and give your answer in symbolic form.

(b) Under what circumstances does your result become infinite? Give a physical interpretation. What happens to your answer when the height of the curb is zero? Does this make sense?

▷ Hint, p. 1035 ■



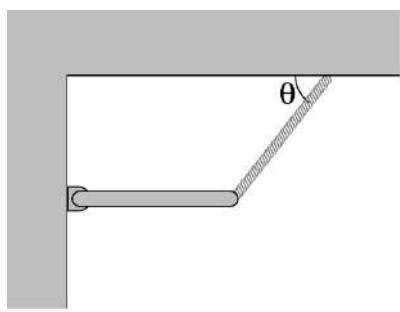
Problem 11.

12 A ball is connected by a string to a vertical post. The ball is set in horizontal motion so that it starts winding the string around the post. Assume that the motion is confined to a horizontal plane, i.e., ignore gravity. Michelle and Astrid are trying to predict the final velocity of the ball when it reaches the post. Michelle says that according to conservation of angular momentum, the ball has to speed up as it approaches the post. Astrid says that according to conservation of energy, the ball has to keep a constant speed. Who is right? [Hint: How is this different from the case where you whirl a rock in a circle on a string and gradually reel in the string?] ■

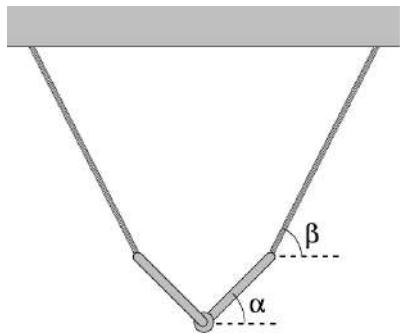
13 In the 1950's, serious articles began appearing in magazines like *Life* predicting that world domination would be achieved by the nation that could put nuclear bombs in orbiting space stations, from which they could be dropped at will. In fact it can be quite difficult to get an orbiting object to come down. Let the object have energy $E = K + U$ and angular momentum L . Assume that the energy is negative, i.e., the object is moving at less than escape velocity. Show that it can never reach a radius less than

$$r_{min} = \frac{GMm}{2E} \left(-1 + \sqrt{1 + \frac{2EL^2}{G^2M^2m^3}} \right).$$

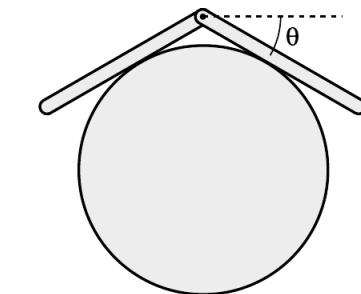
[Note that both factors are negative, giving a positive result.] ■



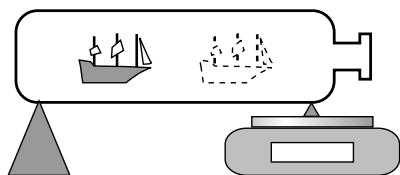
Problem 14.



Problem 15.



Problem 16.



Problem 17.

14 (a) The bar of mass m is attached at the wall with a hinge, and is supported on the right by a massless cable. Find the tension, T , in the cable in terms of the angle θ . ✓

(b) Interpreting your answer to part a, what would be the best angle to use if we wanted to minimize the strain on the cable?

(c) Again interpreting your answer to part a, for what angles does the result misbehave mathematically? Interpret this physically. ■

15 (a) The two identical rods are attached to one another with a hinge, and are supported by the two massless cables. Find the angle α in terms of the angle β , and show that the result is a purely geometric one, independent of the other variables involved. ✓

(b) Using your answer to part a, sketch the configurations for $\beta \rightarrow 0$, $\beta = 45^\circ$, and $\beta = 90^\circ$. Do your results make sense intuitively? ■

16 Two bars of length ℓ are connected with a hinge and placed on a frictionless cylinder of radius r . (a) Show that the angle θ shown in the figure is related to the unitless ratio r/ℓ by the equation

$$\frac{r}{\ell} = \frac{\cos^2 \theta}{2 \tan \theta}.$$

(b) Discuss the physical behavior of this equation for very large and very small values of r/ℓ . ■

17 You wish to determine the mass of a ship in a bottle without taking it out. Show that this can be done with the setup shown in the figure, with a scale supporting the bottle at one end, provided that it is possible to take readings with the ship slid to several different locations. Note that you can't determine the position of the ship's center of mass just by looking at it, and likewise for the bottle. In particular, you can't just say, "position the ship right on top of the fulcrum" or "position it right on top of the balance." ■

18 Suppose that we lived in a universe in which Newton's law of gravity gave an interaction energy proportional to r^{-6} , rather than r^{-1} . Which, if any, of Kepler's laws would still be true? Which would be completely false? Which would be different, but in a way that could be calculated with straightforward algebra? ■

19 Use analogies to find the equivalents of the following equations for rotation in a plane:

$$\begin{aligned} KE &= p^2/2m \\ \Delta x &= v_0\Delta t + (1/2)a\Delta t^2 \\ W &= F\Delta x \end{aligned}$$

Example: $v = \Delta x/\Delta t \rightarrow \omega = \Delta\theta/\Delta t$ ■

20 For a one-dimensional harmonic oscillator, the solution to the energy conservation equation,

$$U + K = \frac{1}{2}kx^2 + \frac{1}{2}mv^2 = \text{constant},$$

is an oscillation with frequency $\omega = \sqrt{k/m}$.

Now consider an analogous system consisting of a bar magnet hung from a thread, which acts like a magnetic compass. A normal compass is full of water, so its oscillations are strongly damped, but the magnet-on-a-thread compass has very little friction, and will oscillate repeatedly around its equilibrium direction. The magnetic energy of the bar magnet is

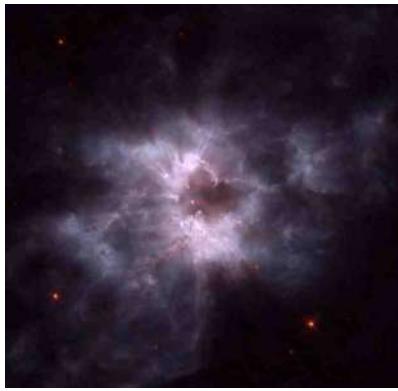
$$U = -Bm \cos \theta,$$

where B is a constant that measures the strength of the earth's magnetic field, m is a constant that parametrizes the strength of the magnet, and θ is the angle, measured in radians, between the bar magnet and magnetic north. The equilibrium occurs at $\theta = 0$, which is the minimum of U .

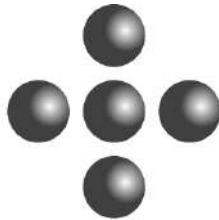
- (a) Problem 19 on p. 297 gave some examples of how to construct analogies between rotational and linear motion. Using the same technique, translate the equation defining the linear quantity k to one that defines an analogous angular one κ (Greek letter kappa). Applying this to the present example, find an expression for κ . (Assume the thread is so thin that its stiffness does not have any significant effect compared to earth's magnetic field.) ✓
(b) Find the frequency of the compass's vibrations. ✓ ■

21 (a) Find the angular velocities of the earth's rotation and of the earth's motion around the sun. ✓

(b) Which motion involves the greater acceleration? ■



Problem 22.



Problem 23

22 The sun turns on its axis once every 26.0 days. Its mass is 2.0×10^{30} kg and its radius is 7.0×10^8 m. Assume it is a rigid sphere of uniform density.

(a) What is the sun's angular momentum? ✓

In a few billion years, astrophysicists predict that the sun will use up all its sources of nuclear energy, and will collapse into a ball of exotic, dense matter known as a white dwarf. Assume that its radius becomes 5.8×10^6 m (similar to the size of the Earth.) Assume it does not lose any mass between now and then. (Don't be fooled by the photo, which makes it look like nearly all of the star was thrown off by the explosion. The visually prominent gas cloud is actually thinner than the best laboratory vacuum ever produced on earth. Certainly a little bit of mass is actually lost, but it is not at all unreasonable to make an approximation of zero loss of mass as we are doing.)

(b) What will its angular momentum be?

(c) How long will it take to turn once on its axis? ✓

23 Give a numerical comparison of the two molecules' moments of inertia for rotation in the plane of the page about their centers of mass. ✓

24 A yo-yo of total mass m consists of two solid cylinders of radius R , connected by a small spindle of negligible mass and radius r . The top of the string is held motionless while the string unrolls from the spindle. Show that the acceleration of the yo-yo is $g/(1 + R^2/2r^2)$. [Hint: The acceleration and the tension in the string are unknown. Use $\tau = \Delta L/\Delta t$ and $F = ma$ to determine these two unknowns.] ■

25 Show that a sphere of radius R that is rolling without slipping has angular momentum and momentum in the ratio $L/p = (2/5)R$. ■

26 Suppose a bowling ball is initially thrown so that it has no angular momentum at all, i.e., it is initially just sliding down the lane. Eventually kinetic friction will get it spinning fast enough so that it is rolling without slipping. Show that the final velocity of the ball equals $5/7$ of its initial velocity. [Hint: You'll need the result of problem 25.] ■

27 Find the angular momentum of a particle whose position is $\mathbf{r} = 3\hat{\mathbf{x}} - \hat{\mathbf{y}} + \hat{\mathbf{z}}$ (in meters) and whose momentum is $\mathbf{p} = -2\hat{\mathbf{x}} + \hat{\mathbf{y}} + \hat{\mathbf{z}}$ (in kg·m/s). ✓ ■

28 Find a vector that is perpendicular to both of the following two vectors:

$$\begin{aligned}\hat{\mathbf{x}} + 2\hat{\mathbf{y}} + 3\hat{\mathbf{z}} \\ 4\hat{\mathbf{x}} + 5\hat{\mathbf{y}} + 6\hat{\mathbf{z}}\end{aligned}$$

✓ ■

29 Prove property (3) of the vector cross product from the theorem on page 1027. ■

30 Prove the anticommutative property of the vector cross product, $\mathbf{A} \times \mathbf{B} = -\mathbf{B} \times \mathbf{A}$, using the expressions for the components of the cross product. Note that giving an example does not constitute a proof of a general rule. ■

31 Find three vectors with which you can demonstrate that the vector cross product need not be associative, i.e., that $\mathbf{A} \times (\mathbf{B} \times \mathbf{C})$ need not be the same as $(\mathbf{A} \times \mathbf{B}) \times \mathbf{C}$. ■

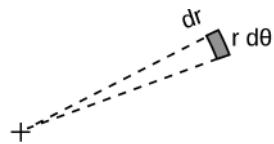
32 Which of the following expressions make sense, and which are nonsense? For those that make sense, indicate whether the result is a vector or a scalar.

- (a) $(\mathbf{A} \times \mathbf{B}) \times \mathbf{C}$
 - (b) $(\mathbf{A} \times \mathbf{B}) \cdot \mathbf{C}$
 - (c) $(\mathbf{A} \cdot \mathbf{B}) \times \mathbf{C}$
-

33 (a) As suggested in the figure, find the area of the infinitesimal region expressed in polar coordinates as lying between r and $r + dr$ and between θ and $\theta + d\theta$. ✓

(b) Generalize this to find the infinitesimal element of volume in cylindrical coordinates (r, θ, z) , where the Cartesian z axis is perpendicular to the directions measured by r and θ . ✓

(c) Find the moment of inertia for rotation about its axis of a cone whose mass is M , whose height is h , and whose base has a radius b . ✓ ■



Problem 33

34 Find the moment of inertia of a solid rectangular box of mass M and uniform density, whose sides are of length a , b , and c , for rotation about an axis through its center parallel to the edges of

length a .

✓ ■

35 The nucleus ^{168}Er (erbium-168) contains 68 protons (which is what makes it a nucleus of the element erbium) and 100 neutrons. It has an ellipsoidal shape like an American football, with one long axis and two short axes that are of equal diameter. Because this is a subatomic system, consisting of only 168 particles, its behavior shows some clear quantum-mechanical properties. It can only have certain energy levels, and it makes quantum leaps between these levels. Also, its angular momentum can only have certain values, which are all multiples of $2.109 \times 10^{-34} \text{ kg} \cdot \text{m}^2/\text{s}$. The table shows some of the observed angular momenta and energies of ^{168}Er , in SI units ($\text{kg} \cdot \text{m}^2/\text{s}$ and joules).

$$L \times 10^{34} \quad E \times 10^{14}$$

0	0
2.109	1.2786
4.218	4.2311
6.327	8.7919
8.437	14.8731
10.546	22.3798
12.655	31.135
14.764	41.206
16.873	52.223

- (a) These data can be described to a good approximation as a rigid end-over-end rotation. Estimate a single best-fit value for the moment of inertia from the data, and check how well the data agree with the assumption of rigid-body rotation. \triangleright Hint, p. 1035 \checkmark
 (b) Check whether this moment of inertia is on the right order of magnitude. The moment of inertia depends on both the size and the shape of the nucleus. For the sake of this rough check, ignore the fact that the nucleus is not quite spherical. To estimate its size, use the fact that a neutron or proton has a volume of about 1 fm^3 (one cubic femtometer, where $1 \text{ fm} = 10^{-15} \text{ m}$), and assume they are closely packed in the nucleus. \blacksquare

36 (a) Prove the identity $\mathbf{a} \times (\mathbf{b} \times \mathbf{c}) = \mathbf{b}(\mathbf{a} \cdot \mathbf{c}) - \mathbf{c}(\mathbf{a} \cdot \mathbf{b})$ by expanding the product in terms of its components. Note that because the x , y , and z components are treated symmetrically in the definitions of the vector cross product, it is only necessary to carry out the proof for the x component of the result.

(b) Applying this to the angular momentum of a rigidly rotating body, $L = \int \mathbf{r} \times (\boldsymbol{\omega} \times \mathbf{r}) dm$, show that the diagonal elements of the moment of inertia tensor can be expressed as, e.g., $I_{xx} = \int (y^2 + z^2) dm$.

(c) Find the diagonal elements of the moment of inertia matrix of an ellipsoid with axes of lengths a , b , and c , in the principal-axis frame, and with the axis at the center. \checkmark \blacksquare

37 In example 22 on page 282, prove that if the rod is sufficiently thin, it can be toppled without scraping on the floor.

\triangleright Solution, p. 1044 \triangleright Solution, p. 1044 \blacksquare

38 Suppose an object has mass m , and moment of inertia I_0 for rotation about some axis A passing through its center of mass. Prove that for an axis B, parallel to A and lying at a distance h from it, the object's moment of inertia is given by $I_0 + mh^2$. This is known as the parallel axis theorem. ■

39 Let two sides of a triangle be given by the vectors \mathbf{A} and \mathbf{B} , with their tails at the origin, and let mass m be uniformly distributed on the interior of the triangle. (a) Show that the distance of the triangle's center of mass from the intersection of sides \mathbf{A} and \mathbf{B} is given by $\frac{1}{3}|\mathbf{A} + \mathbf{B}|$.

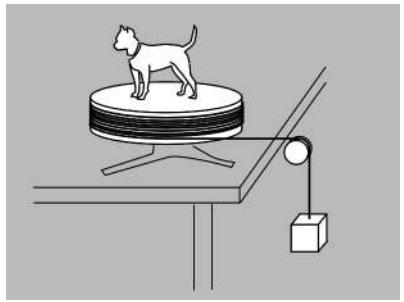
(b) Consider the quadrilateral with mass $2m$, and vertices at the origin, \mathbf{A} , \mathbf{B} , and $\mathbf{A} + \mathbf{B}$. Show that its moment of inertia, for rotation about an axis perpendicular to it and passing through its center of mass, is $\frac{m}{6}(A^2 + B^2)$.

(c) Show that the moment of inertia for rotation about an axis perpendicular to the plane of the original triangle, and passing through its center of mass, is $\frac{m}{18}(A^2 + B^2 - \mathbf{A} \cdot \mathbf{B})$. Hint: Combine the results of parts a and b with the result of problem 38.

40 When we talk about rigid-body rotation, the concept of a perfectly rigid body can only be an idealization. In reality, any object will compress, expand, or deform to some extent when subjected to the strain of rotation. However, if we let it settle down for a while, perhaps it will reach a new equilibrium. As an example, suppose we fill a centrifuge tube with some compressible substance like shaving cream or Wonder Bread. We can model the contents of the tube as a one-dimensional line of mass, extending from $r = 0$ to $r = \ell$. Once the rotation starts, we expect that the contents will be most compressed near the “floor” of the tube at $r = \ell$; this is both because the inward force required for circular motion increases with r for a fixed ω , and because the part at the floor has the greatest amount of material pressing “down” (actually outward) on it. The linear density dm/dr , in units of kg/m, should therefore increase as a function of r . Suppose that we have $dm/dr = \mu e^{r/\ell}$, where μ is a constant. Find the moment of inertia. ✓ ■

41 When we release an object such as a bicycle wheel or a coin on an inclined plane, we can observe a variety of different behaviors. Characterize these behaviors empirically and try to list the physical parameters that determine which behavior occurs. Try to form a conjecture about the behavior using simple closed-form expressions. Test your conjecture experimentally. ■

42 The figure shows a tabletop experiment that can be used to determine an unknown moment of inertia. A rotating platform of radius R has a string wrapped around it. The string is threaded over a pulley and down to a hanging weight of mass m . The mass is released from rest, and its downward acceleration a ($a > 0$) is measured. Find the total moment of inertia I of the platform plus the object sitting on top of it. (The moment of inertia of the object itself can then be found by subtracting the value for the empty platform.) \checkmark



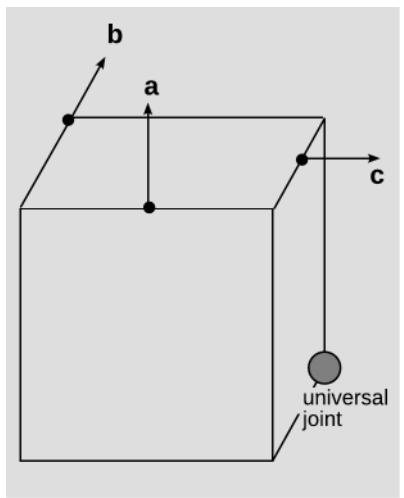
Problem 42.

43 The uniform cube has unit weight and sides of unit length. One corner is attached to a universal joint, i.e., a frictionless bearing that allows any type of rotation. If the cube is in equilibrium, find the magnitudes of the forces \mathbf{a} , \mathbf{b} , and \mathbf{c} . \checkmark

44 In this problem we investigate the notion of division by a vector.

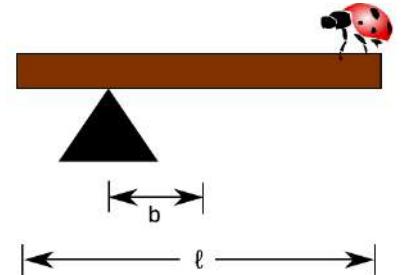
- (a) Given a nonzero vector \mathbf{a} and a scalar b , suppose we wish to find a vector \mathbf{u} that is the solution of $\mathbf{a} \cdot \mathbf{u} = b$. Show that the solution is not unique, and give a geometrical description of the solution set.
- (b) Do the same thing for the equation $\mathbf{a} \times \mathbf{u} = \mathbf{c}$.
- (c) Show that the *simultaneous* solution of these two equations exists and is unique.

Remark: This is one motivation for constructing the number system called the quaternions. For a certain period around 1900, quaternions were more popular than the system of vectors and scalars more commonly used today. They still have some important advantages over the scalar-vector system for certain applications, such as avoiding a phenomenon known as gimbal lock in controlling the orientation of bodies such as spacecraft. \blacksquare



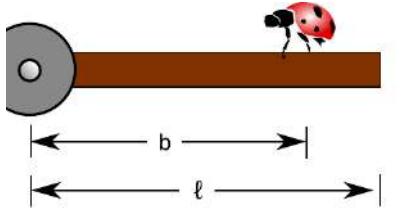
Problem 43.

45 Show that when a thin, uniform ring rotates about a diameter, the moment of inertia is half as big as for rotation about the axis of symmetry. \triangleright Solution, p. 1044 \blacksquare



Problem 46.

46 A bug stands at the right end of a rod of length ℓ , which is initially at rest in a horizontal position. The rod rests on a fulcrum which is at a distance b to the left of the rod's center, so that when the rod is released from rest, the bug's end will drop. For what value of b will the bug experience apparent weightlessness at the moment when the rod is released? \checkmark



Problem 47.

47 The figure shows a trap door of length ℓ , which is released at rest from a horizontal position and swings downward under its own weight. The bug stands at a distance b from the hinge. Because the bug feels the floor dropping out from under it with some acceleration, it feels a change in the apparent acceleration of gravity from g to some value g_a , at the moment when the door is released. Find g_a . \checkmark

Key to symbols:

\blacksquare easy \blacksquare typical \blacksquare challenging \blacksquare difficult \blacksquare very difficult

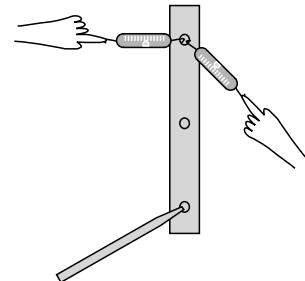
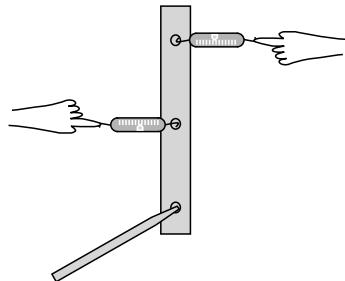
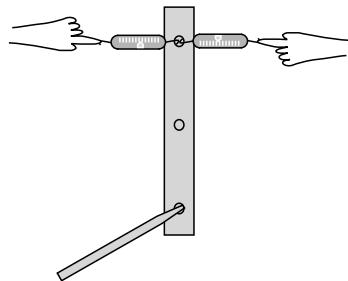
✓ An answer check is available at www.lightandmatter.com.

Exercises

Exercise 4A: Torque

Equipment:

- rulers with holes in them
- spring scales (two per group)



While one person holds the pencil which forms the axle for the ruler, the other members of the group pull on the scale and take readings. In each case, calculate the total torque on the ruler, and find out whether it equals zero to roughly within the accuracy of the experiment.

Chapter 5

Thermodynamics

$$S = k \log W$$

*Inscription on the tomb of Ludwig Boltzmann, 1844-1906.
Boltzmann originated the microscopic theory of thermodynamics.*

In a developing country like China, a refrigerator is the mark of a family that has arrived in the middle class, and a car is the ultimate symbol of wealth. Both of these are *heat engines*: devices for converting between heat and other forms of energy. Unfortunately for the Chinese, neither is a very efficient device. Burning fossil fuels has made China's big cities the most polluted on the planet, and the country's total energy supply isn't sufficient to support American levels of energy consumption by more than a small fraction of China's population. Could we somehow manipulate energy in a more efficient way?

Conservation of energy is a statement that the total amount of energy is constant at all times, which encourages us to believe that any energy transformation can be undone — indeed, the laws of physics you've learned so far don't even distinguish the past from the future. If you get in a car and drive around the block, the net effect is to consume some of the energy you paid for at the gas station, using it to heat the neighborhood. There would not seem to be any fundamental physical principle to prevent you from recapturing all that heat and using it again the next time you want to go for a drive. More modestly, why don't engineers design a car engine so that it recaptures the heat energy that would otherwise be wasted via the radiator and the exhaust?

Hard experience, however, has shown that designers of more and more efficient engines run into a brick wall at a certain point. The generators that the electric company uses to produce energy at an oil-fueled plant are indeed much more efficient than a car engine, but even if one is willing to accept a device that is very large, expensive, and complex, it turns out to be impossible to make a perfectly efficient heat engine — not just impossible with present-day technology, but impossible due to a set of fundamental physical principles known as the science of *thermodynamics*. And thermodynamics isn't just a pesky set of constraints on heat engines. Without thermodynamics, there is no way to explain the direction of time's arrow — why we

can remember the past but not the future, and why it's easier to break Humpty Dumpty than to put him back together again.

5.1 Pressure, temperature, and heat

When we heat an object, we speed up the mind-bogglingly complex random motion of its molecules. One method for taming complexity is the conservation laws, since they tell us that certain things must remain constant regardless of what process is going on. Indeed, the law of conservation of energy is also known as the first law of thermodynamics.

But as alluded to in the introduction to this chapter, conservation of energy by itself is not powerful enough to explain certain empirical facts about heat. A second way to sidestep the complexity of heat is to ignore heat's atomic nature and concentrate on quantities like temperature and pressure that tell us about a system's properties as a whole. This approach is called macroscopic in contrast to the microscopic method of attack. Pressure and temperature were fairly well understood in the age of Newton and Galileo, hundreds of years before there was any firm evidence that atoms and molecules even existed.

Unlike the conserved quantities such as mass, energy, momentum, and angular momentum, neither pressure nor temperature is additive. Two cups of coffee have twice the heat energy of a single cup, but they do not have twice the temperature. Likewise, the painful pressure on your eardrums at the bottom of a pool is not affected if you insert or remove a partition between the two halves of the pool.

We restrict ourselves to a discussion of pressure in fluids at rest and in equilibrium. In physics, the term "fluid" is used to mean either a gas or a liquid. The important feature of a fluid can be demonstrated by comparing with a cube of jello on a plate. The jello is a solid. If you shake the plate from side to side, the jello will respond by shearing, i.e., by slanting its sides, but it will tend to spring back into its original shape. A solid can sustain shear forces, but a fluid cannot. A fluid does not resist a change in shape unless it involves a change in volume.

5.1.1 Pressure

If you're at the bottom of a pool, you can't relieve the pain in your ears by turning your head. The water's force on your eardrum is always the same, and is always perpendicular to the surface where the eardrum contacts the water. If your ear is on the east side of your head, the water's force is to the west. If you keep your ear in the same spot while turning around so your ear is on the north, the force will still be the same in magnitude, and it will change its direction so that it is still perpendicular to the eardrum: south.

This shows that pressure has no direction in space, i.e., it is a scalar. The direction of the force is determined by the orientation of the surface on which the pressure acts, not by the pressure itself. A fluid flowing over a surface can also exert frictional forces, which are parallel to the surface, but the present discussion is restricted to fluids at rest.

Experiments also show that a fluid's force on a surface is proportional to the surface area. The vast force of the water behind a dam, for example, in proportion to the dam's great surface area. (The bottom of the dam experiences a higher proportion of its force.)

Based on these experimental results, it appears that the useful way to define pressure is as follows. The pressure of a fluid at a given point is defined as F_{\perp}/A , where A is the area of a small surface inserted in the fluid at that point, and F_{\perp} is the component of the fluid's force on the surface which is perpendicular to the surface. (In the case of a moving fluid, fluid friction forces can act parallel to the surface, but we're only dealing with stationary fluids, so there is only an F_{\perp} .)

This is essentially how a pressure gauge works. The reason that the surface must be small is so that there will not be any significant difference in pressure between one part of it and another part. The SI units of pressure are evidently N/m^2 , and this combination can be abbreviated as the pascal, $1 \text{ Pa}=1 \text{ N/m}^2$. The pascal turns out to be an inconveniently small unit, so car tires, for example, normally have pressures imprinted on them in units of kilopascals.

Pressure in U.S. units

example 1

In U.S. units, the unit of force is the pound, and the unit of distance is the inch. The unit of pressure is therefore pounds per square inch, or p.s.i. (Note that the pound is not a unit of mass.)

Atmospheric pressure in U.S. and metric units

example 2

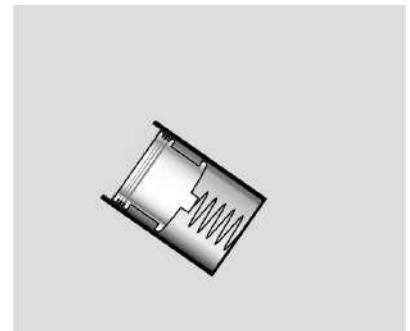
▷ A figure that many people in the U.S. remember is that atmospheric pressure is about 15 pounds per square inch. What is this in metric units?

▷

$$\begin{aligned}(15 \text{ lb})/(1 \text{ in}^2) &= \frac{68 \text{ N}}{(0.0254 \text{ m})^2} \\ &= 1.0 \times 10^5 \text{ N/m}^2 \\ &= 100 \text{ kPa}\end{aligned}$$

Only pressure differences are normally significant.

If you spend enough time on an airplane, the pain in your ears subsides. This is because your body has gradually been able to admit more air into the cavity behind the eardrum. Once the pressure



a / A simple pressure gauge consists of a cylinder open at one end, with a piston and a spring inside. The depth to which the spring is depressed is a measure of the pressure. To determine the absolute pressure, the air needs to be pumped out of the interior of the gauge, so that there is no air pressure acting outward on the piston. In many practical gauges, the back of the piston is open to the atmosphere, so the pressure the gauge registers equals the pressure of the fluid minus the pressure of the atmosphere.

inside is equalized with the pressure outside, the inward and outward forces on your eardrums cancel out, and there is no physical sensation to tell you that anything unusual is going on. For this reason, it is normally only pressure differences that have any physical significance. Thus deep-sea fish are perfectly healthy in their habitat because their bodies have enough internal pressure to cancel the pressure from the water in which they live; if they are caught in a net and brought to the surface rapidly, they explode because their internal pressure is so much greater than the low pressure outside.

Getting killed by a pool pump

example 3

▷ My house has a pool, which I maintain myself. A pool always needs to have its water circulated through a filter for several hours a day in order to keep it clean. The filter is a large barrel with a strong clamp that holds the top and bottom halves together. My filter has a prominent warning label that warns me not to try to open the clamps while the pump is on, and it shows a cartoon of a person being struck by the top half of the pump. The cross-sectional area of the filter barrel is 0.25 m^2 . Like most pressure gauges, the one on my pool pump actually reads the difference in pressure between the pressure inside the pump and atmospheric pressure. The gauge reads 90 kPa. What is the force that is trying to pop open the filter?

▷ If the gauge told us the absolute pressure of the water inside, we'd have to find the force of the water pushing outward and the force of the air pushing inward, and subtract in order to find the total force. Since air surrounds us all the time, we would have to do such a subtraction every time we wanted to calculate anything useful based on the gauge's reading. The manufacturers of the gauge decided to save us from all this work by making it read the difference in pressure between inside and outside, so all we have to do is multiply the gauge reading by the cross-sectional area of the filter:

$$\begin{aligned} F &= PA \\ &= (90 \times 10^3 \text{ N/m}^2)(0.25 \text{ m}^2) \\ &= 22000 \text{ N} \end{aligned}$$

That's a lot of force!

The word “suction” and other related words contain a hidden misunderstanding related to this point about pressure differences. When you suck water up through a straw, there is nothing in your mouth that is attracting the water upward. The force that lifts the water is from the pressure of the water in the cup. By creating a partial vacuum in your mouth, you decreased the air's downward force on the water so that it no longer exactly canceled the upward force.

Variation of pressure with depth

The pressure within a fluid in equilibrium can only depend on depth, due to gravity. If the pressure could vary from side to side, then a piece of the fluid in between, b, would be subject to unequal forces from the parts of the fluid on its two sides. Since fluids do not exhibit shear forces, there would be no other force that could keep this piece of fluid from accelerating. This contradicts the assumption that the fluid was in equilibrium.

self-check A

How does this proof fail for solids?

▷ Answer, p. 1061

To find the variation with depth, we consider the vertical forces acting on a tiny, imaginary cube of the fluid having infinitesimal height dy and areas dA on the top and bottom. Using positive numbers for upward forces, we have

$$P_{bottom} dA - P_{top} dA - F_g = 0.$$

The weight of the fluid is $F_g = mg = \rho V g = \rho dA dy g$, where ρ is the density of the fluid, so the difference in pressure is

$$dP = -\rho g dy. \quad \begin{aligned} & \text{[variation in pressure with depth for} \\ & \text{a fluid of density } \rho \text{ in equilibrium;} \\ & \text{positive } y \text{ is up.]} \end{aligned}$$

A more elegant way of writing this is in terms of a dot product, $dP = \rho g \cdot dy$, which automatically takes care of the plus or minus sign, depending on the relative directions of the g and dy vectors, and avoids any requirements about the coordinate system.

The factor of ρ explains why we notice the difference in pressure when diving 3 m down in a pool, but not when going down 3 m of stairs. The equation only tells us the difference in pressure, not the absolute pressure. The pressure at the surface of a swimming pool equals the atmospheric pressure, not zero, even though the depth is zero at the surface. The blood in your body does not even have an upper surface.

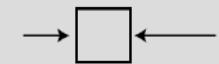
In cases where g and ρ are independent of depth, we can integrate both sides of the equation to get everything in terms of finite differences rather than differentials: $\Delta P = -\rho g \Delta y$.

self-check B

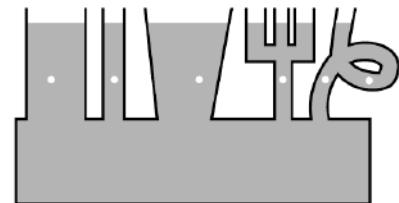
In which of the following situations is the equation $\Delta P = -\rho g \Delta y$ valid?

Why? (1) difference in pressure between a tabletop and the feet (i.e., predicting the pressure of the feet on the floor) (2) difference in air pressure between the top and bottom of a tall building (3) difference in air pressure between the top and bottom of Mt. Everest (4) difference in pressure between the top of the earth's mantle and the center of the earth (5) difference in pressure between the top and bottom of an airplane's wing

▷ Answer, p. 1061



b / This doesn't happen. If pressure could vary horizontally in equilibrium, the cube of water would accelerate horizontally. This is a contradiction, since we assumed the fluid was in equilibrium.



c / The pressure is the same at all the points marked with dots.



d / This does happen. The sum of the forces from the surrounding parts of the fluid is upward, canceling the downward force of gravity.

Pressure of lava underneath a volcano

example 4

▷ A volcano has just finished erupting, and a pool of molten lava is lying at rest in the crater. The lava has come up through an opening inside the volcano that connects to the earth's molten mantle. The density of the lava is 4.1 g/cm^3 . What is the pressure in the lava underneath the base of the volcano, 3000 m below the surface of the pool?

▷

$$\begin{aligned}\Delta P &= \rho g \Delta y \\ &= (4.1 \times 10^3 \text{ kg/m}^3)(9.8 \text{ m/s}^2)(3000 \text{ m}) \\ &= 1.2 \times 10^8 \text{ Pa}\end{aligned}$$

This is the difference between the pressure we want to find and atmospheric pressure at the surface. The latter, however, is tiny compared to the ΔP we just calculated, so what we've found is essentially the pressure, P .

Atmospheric pressure

example 5

Gases, unlike liquids, are quite compressible, and at a given temperature, the density of a gas is approximately proportional to the pressure. The proportionality constant is discussed on page 317, but for now let's just call it k , $\rho = kP$. Using this fact, we can find the variation of atmospheric pressure with altitude, assuming constant temperature:

$$dP = -\rho g dy = -kPg dy$$

$$\frac{dP}{P} = -kg dy$$

$$\ln P = -kg y + \text{constant} \quad [\text{integrating both sides}]$$

$$P = (\text{constant})e^{-kg y} \quad [\text{exponentiating both sides}]$$



e / We have to wait for the thermometer to equilibrate its temperature with the temperature of Irene's armpit.

5.1.2 Temperature

Thermal equilibrium

We use the term temperature casually, but what is it exactly? Roughly speaking, temperature is a measure of how concentrated the heat energy is in an object. A large, massive object with very little heat energy in it has a low temperature.

But physics deals with operational definitions, i.e., definitions of how to measure the thing in question. How do we measure temperature? One common feature of all temperature-measuring devices is that they must be left for a while in contact with the thing whose temperature is being measured. When you take your temperature

with a fever thermometer, you are waiting for the mercury inside to come up to the same temperature as your body. The thermometer actually tells you the temperature of its own working fluid (in this case the mercury). In general, the idea of temperature depends on the concept of thermal equilibrium. When you mix cold eggs from the refrigerator with flour that has been at room temperature, they rapidly reach a compromise temperature. What determines this compromise temperature is conservation of energy, and the amount of energy required to heat or cool each substance by one degree. But without even having constructed a temperature scale, we can see that the important point is the phenomenon of thermal equilibrium itself: two objects left in contact will approach the same temperature. We also assume that if object A is at the same temperature as object B, and B is at the same temperature as C, then A is at the same temperature as C. This statement is sometimes known as the zeroth law of thermodynamics, so called because after the first, second, and third laws had been developed, it was realized that there was another law that was even more fundamental.

Thermal expansion

The familiar mercury thermometer operates on the principle that the mercury, its working fluid, expands when heated and contracts when cooled. In general, all substances expand and contract with changes in temperature. The zeroth law of thermodynamics guarantees that we can construct a comparative scale of temperatures that is independent of what type of thermometer we use. If a thermometer gives a certain reading when it's in thermal equilibrium with object A, and also gives the same reading for object B, then A and B must be the same temperature, regardless of the details of how the thermometers works.

What about constructing a temperature scale in which every degree represents an equal step in temperature? The Celsius scale has 0 as the freezing point of water and 100 as its boiling point. The hidden assumption behind all this is that since two points define a line, any two thermometers that agree at two points must agree at all other points. In reality if we calibrate a mercury thermometer and an alcohol thermometer in this way, we will find that a graph of one thermometer's reading versus the other is not a perfectly straight $y = x$ line. The subtle inconsistency becomes a drastic one when we try to extend the temperature scale through the points where mercury and alcohol boil or freeze. Gases, however, are much more consistent among themselves in their thermal expansion than solids or liquids, and the noble gases like helium and neon are more consistent with each other than gases in general. Continuing to search for consistency, we find that noble gases are more consistent with each other when their pressure is very low.

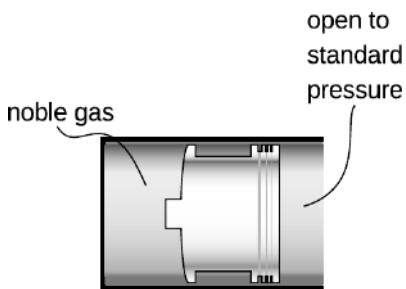
As an idealization, we imagine a gas in which the atoms interact



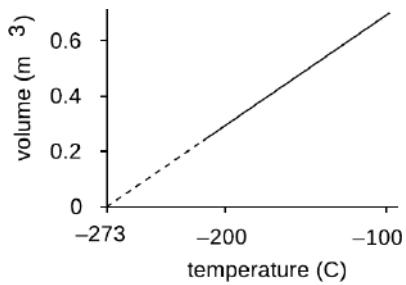
f / Thermal equilibrium can be prevented. Otters have a coat of fur that traps air bubbles for insulation. If a swimming otter was in thermal equilibrium with cold water, it would be dead. Heat is still conducted from the otter's body to the water, but much more slowly than it would be in a warm-blooded animal that didn't have this special adaptation.



g / A hot air balloon is inflated. Because of thermal expansion, the hot air is less dense than the surrounding cold air, and therefore floats as the cold air drops underneath it and pushes it up out of the way.



h / A simplified version of an ideal gas thermometer. The whole instrument is allowed to come into thermal equilibrium with the substance whose temperature is to be measured, and the mouth of the cylinder is left open to standard pressure. The volume of the noble gas gives an indication of temperature.



i / The volume of 1 kg of neon gas as a function of temperature (at standard pressure). Although neon would actually condense into a liquid at some point, extrapolating the graph gives to zero volume gives the same temperature as for any other gas: absolute zero.

only with the sides of the container, not with each other. Such a gas is perfectly nonreactive (as the noble gases very nearly are), and never condenses to a liquid (as the noble gases do only at extremely low temperatures). Its atoms take up a negligible fraction of the available volume. Any gas can be made to behave very much like this if the pressure is extremely low, so that the atoms hardly ever encounter each other. Such a gas is called an ideal gas, and we define the Celsius scale in terms of the volume of the gas in a thermometer whose working substance is an ideal gas maintained at a fixed (very low) pressure, and which is calibrated at 0 and 100 degrees according to the melting and boiling points of water. The Celsius scale is not just a comparative scale but an additive one as well: every step in temperature is equal, and it makes sense to say that the difference in temperature between 18 and 28°C is the same as the difference between 48 and 58.

Absolute zero and the kelvin scale

We find that if we extrapolate a graph of volume versus temperature, the volume becomes zero at nearly the same temperature for all gases: -273°C . Real gases will all condense into liquids at some temperature above this, but an ideal gas would achieve zero volume at this temperature, known as absolute zero. The most useful temperature scale in scientific work is one whose zero is defined by absolute zero, rather than by some arbitrary standard like the melting point of water. The temperature scale used universally in scientific work, called the Kelvin scale, is the same as the Celsius scale, but shifted by 273 degrees to make its zero coincide with absolute zero. Scientists use the Celsius scale only for comparisons or when a change in temperature is all that is required for a calculation. Only on the Kelvin scale does it make sense to discuss ratios of temperatures, e.g., to say that one temperature is twice as hot as another.

Which temperature scale to use

example 6

- ▷ You open an astronomy book and encounter the equation

$$(\text{light emitted}) = (\text{constant}) \times T^4$$

for the light emitted by a star as a function of its surface temperature. What temperature scale is implied?

- ▷ The equation tells us that doubling the temperature results in the emission of 16 times as much light. Such a ratio only makes sense if the Kelvin scale is used.

Although we can achieve as good an approximation to an ideal gas as we wish by making the pressure very low, it seems nevertheless that there should be some more fundamental way to define temperature. We will construct a more fundamental scale of temperature in section 5.4.

5.1.3 Heat

“Heat,” notated Q , is used in thermodynamics as a term for an amount of thermal energy that is transferred. When you put a bite of food in your mouth that is too hot, the pain is caused by the heat transferred from the food to your mouth. People discussing the weather may say “What about this heat today?” or “What about this temperature today?” as if the words were synonyms, but to a physicist they are distinct. Temperature is not additive, but heat is: two sips of hot coffee have the same temperature as one, but two sips will transfer twice the heat to your mouth. Temperature is measured in degrees, heat in joules.

If I give you an object, you can measure its temperature — physicists call temperature a “property of state,” i.e., you can tell what it is from the current state of the object. Heat is a description of a *process* of energy transfer, not a property of state.

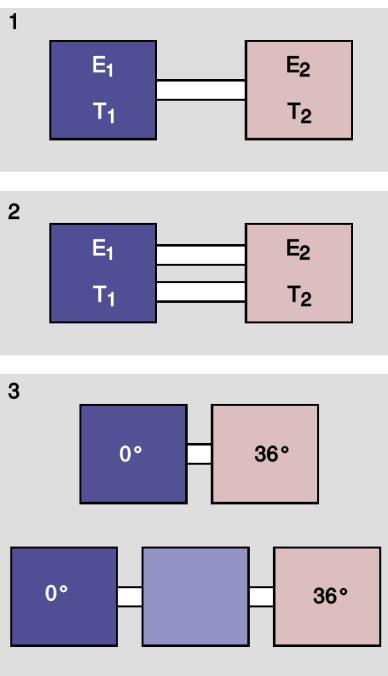
It’s relatively easy to detect and measure a *transfer* of thermal energy (the hot bite of food), but to say how much thermal energy an object *has* is much harder — sometimes even impossible in principle.

Heat is distinguished from mechanical work (3.2.8, p. 164) because work is the transfer of energy by a macroscopically measurable force, e.g., the force of a baseball bat on the ball. No such force is needed in order to melt an ice cube; the forces are in microscopic collisions of water molecules with ice molecules.

Heat, like the flow of money or water, is a signed quantity, but the sign is a matter of definition. The bank’s debit is the customer’s withdrawal. It is an arbitrary choice whether to call Q positive when it flows from object A to object B or from B to A, and likewise for the work W . Similar choices arise in the description of flowing fluids (example 1, p. 58) or electric currents (sec. 9.1.1, p. 532). We will usually adopt definitions such that as many heats and works as possible are positive. So by our definition, a cute 19th-century steam locomotive takes in positive heat from its boiler, does positive work to pull the cars, and spews out positive heat through its smokestack. When only a single object is being discussed, such as a cylinder of compressed air, we define a heat input as positive and a work output as positive, which is again in accord with the picture of the cute steam engine. No universally consistent convention is possible, since, e.g., if objects A, B, and C all interact, we will always have opposite signs for A’s work on B and B’s work on A, etc.

Discussion Questions

A Figure j/1 shows objects 1 and 2, each with a certain temperature T and a certain amount of thermal energy E . They are connected by a thin rod, so that eventually they will reach thermal equilibrium. We expect that the rate at which heat is transferred into object 1 will be given by some equation $dE_1/dt = k(\dots)$, where k is a positive constant of proportionality and “ \dots ” is some expression that depends on the temperatures. Suppose



j / Discussion questions A-C.

that the following six forms are proposed for the “...” in $dE_1/dt = k(\dots)$.

1. T_1
2. T_2
3. $T_1 - T_2$
4. $T_2 - T_1$
5. T_1/T_2
6. T_2/T_1

Give physical reasons why five of these are not possible.

B How should the rate of heat conduction in j/2 compare with the rate in j/1?

C The example in j/3 is different from the preceding ones because when we add the third object in the middle, we don't necessarily know the intermediate temperature. We could in fact set up this third object with any desired initial temperature. Suppose, however, that the flow of heat is *steady*. For example, the 36° object could be a human body, the 0° object could be the air on a cold day, and the object in between could be a simplified physical model of the insulation provided by clothing or body fat. Under this assumption, what is the intermediate temperature? How does the rate of heat conduction compare in the two cases?

D Based on the conclusions of questions A-C, how should the rate of heat conduction through an object depend on its length and cross-sectional area? If all the linear dimensions of the object are doubled, what happens to the rate of heat conduction through it? How would this apply if we compare an elephant to a shrew?

5.2 Microscopic description of an ideal gas

5.2.1 Evidence for the kinetic theory

Why does matter have the thermal properties it does? The basic answer must come from the fact that matter is made of atoms. How, then, do the atoms give rise to the bulk properties we observe? Gases, whose thermal properties are simple, offer the best opportunity for us to construct a simple connection between the microscopic and macroscopic worlds.

A crucial observation is that although solids and liquids are nearly incompressible, gases can be compressed, as when we increase the amount of air in a car's tire while hardly increasing its volume at all. This makes us suspect that the atoms in a solid are packed shoulder to shoulder, while a gas is mostly vacuum, with large spaces between molecules. Most liquids and solids have densities about 1000 times greater than most gases, so evidently each molecule in a gas is separated from its nearest neighbors by a space something like 10 times the size of the molecules themselves.

If gas molecules have nothing but empty space between them, why don't the molecules in the room around you just fall to the

floor? The only possible answer is that they are in rapid motion, continually rebounding from the walls, floor and ceiling. In section 2.4 I have already given some of the evidence for the kinetic theory of heat, which states that heat is the kinetic energy of randomly moving molecules. This theory was proposed by Daniel Bernoulli in 1738, and met with considerable opposition because it seemed as though the molecules in a gas would eventually calm down and settle into a thin film on the floor. There was no precedent for this kind of perpetual motion. No rubber ball, however elastic, rebounds from a wall with exactly as much energy as it originally had, nor do we ever observe a collision between balls in which none of the kinetic energy at all is converted to heat and sound. The analogy is a false one, however. A rubber ball consists of atoms, and when it is heated in a collision, the heat is a form of motion of those atoms. An individual molecule, however, cannot possess heat. Likewise sound is a form of bulk motion of molecules, so colliding molecules in a gas cannot convert their kinetic energy to sound. Molecules can indeed induce vibrations such as sound waves when they strike the walls of a container, but the vibrations of the walls are just as likely to impart energy to a gas molecule as to take energy from it. Indeed, this kind of exchange of energy is the mechanism by which the temperatures of the gas and its container become equilibrated.

5.2.2 Pressure, volume, and temperature

A gas exerts pressure on the walls of its container, and in the kinetic theory we interpret this apparently constant pressure as the averaged-out result of vast numbers of collisions occurring every second between the gas molecules and the walls. The empirical facts about gases can be summarized by the relation

$$PV \propto nT, \quad [\text{ideal gas}]$$

which really only holds exactly for an ideal gas. Here n is the number of molecules in the sample of gas.

Volume related to temperature *example 7*

The proportionality of volume to temperature at fixed pressure was the basis for our definition of temperature.

Pressure related to temperature *example 8*

Pressure is proportional to temperature when volume is held constant. An example is the increase in pressure in a car's tires when the car has been driven on the freeway for a while and the tires and air have become hot.

We now connect these empirical facts to the kinetic theory of a classical ideal gas. For simplicity, we assume that the gas is monoatomic (i.e., each molecule has only one atom), and that it is confined to a cubical box of volume V , with L being the length of each edge and A the area of any wall. An atom whose velocity has an x component v_x will collide regularly with the left-hand wall,

traveling a distance $2L$ parallel to the x axis between collisions with that wall. The time between collisions is $\Delta t = 2L/v_x$, and in each collision the x component of the atom's momentum is reversed from $-mv_x$ to mv_x . The total force on the wall is

$$F = \sum \frac{\Delta p_{x,i}}{\Delta t_i} \quad [\text{monoatomic ideal gas}],$$

where the index i refers to the individual atoms. Substituting $\Delta p_{x,i} = 2mv_{x,i}$ and $\Delta t_i = 2L/v_{x,i}$, we have

$$F = \frac{1}{L} \sum mv_{x,i}^2 \quad [\text{monoatomic ideal gas}].$$

The quantity $mv_{x,i}^2$ is twice the contribution to the kinetic energy from the part of the atoms' center of mass motion that is parallel to the x axis. Since we're assuming a monoatomic gas, center of mass motion is the only type of motion that gives rise to kinetic energy. (A more complex molecule could rotate and vibrate as well.) If the quantity inside the sum included the y and z components, the sum would be twice the total kinetic energy of all the molecules. Since we expect the energy to be equally shared among x , y , and z motion,¹ the quantity inside the sum must therefore equal $2/3$ of the total kinetic energy, so

$$F = \frac{2K_{total}}{3L} \quad [\text{monoatomic ideal gas}].$$

Dividing by A and using $AL = V$, we have

$$P = \frac{2K_{total}}{3V} \quad [\text{monoatomic ideal gas}].$$

This can be connected to the empirical relation $PV \propto nT$ if we multiply by V on both sides and rewrite K_{total} as $n\bar{K}$, where \bar{K} is the average kinetic energy per molecule:

$$PV = \frac{2}{3}n\bar{K} \quad [\text{monoatomic ideal gas}].$$

For the first time we have an interpretation of temperature based on a microscopic description of matter: in a monoatomic ideal gas, the temperature is a measure of the average kinetic energy per molecule. The proportionality between the two is $\bar{K} = (3/2)kT$, where the constant of proportionality k , known as Boltzmann's constant, has a numerical value of 1.38×10^{-23} J/K.

The Boltzmann constant has the value it does because the Celsius and kelvin scales were defined before the microscopic picture of thermodynamics had been discovered. For some calculations, it is more convenient to work in more natural units where $k = 1$ by definition, and then the units of temperature and energy are the

¹This equal sharing will be justified more rigorously on page 333.

same. The Boltzmann constant is small because our energy scale of joules is a macroscopic scale, so that when we express the thermal energy of a single atom in joules, the number is very small.

Summarizing, we have the following two important facts.

Microscopic model of an ideal gas

For an ideal gas,

$$PV = nkT,$$

which is known as the ideal gas law. The temperature of the gas is a measure of the average kinetic energy per atom,

$$\bar{K} = \frac{3}{2}kT.$$

Although I won't prove it here, the ideal gas law applies to all ideal gases, even though the derivation assumed a monoatomic ideal gas in a cubical box. (You may have seen it written elsewhere as $PV = NRT$, where $N = n/N_A$ is the number of moles of atoms, $R = kN_A$, and $N_A = 6.0 \times 10^{23}$, called Avogadro's number, is essentially the number of hydrogen atoms in 1 g of hydrogen.)

Pressure in a car tire

example 9

- ▷ After driving on the freeway for a while, the air in your car's tires heats up from 10°C to 35°C . How much does the pressure increase?
- ▷ The tires may expand a little, but we assume this effect is small, so the volume is nearly constant. From the ideal gas law, the ratio of the pressures is the same as the ratio of the absolute temperatures,

$$\begin{aligned} P_2/P_1 &= T_2/T_1 \\ &= (308\text{ K})/(283\text{ K}) \\ &= 1.09, \end{aligned}$$

or a 9% increase.

Discussion Questions

- A** Compare the amount of energy needed to heat 1 liter of helium by 1 degree with the energy needed to heat 1 liter of xenon. In both cases, the heating is carried out in a sealed vessel that doesn't allow the gas to expand. (The vessel is also well insulated.)
- B** Repeat discussion question A if the comparison is 1 kg of helium versus 1 kg of xenon (equal masses, rather than equal volumes).
- C** Repeat discussion question A, but now compare 1 liter of helium in a vessel of constant volume with the same amount of helium in a vessel that allows expansion beyond the initial volume of 1 liter. (This could be a piston, or a balloon.)

5.3 Entropy as a macroscopic quantity

5.3.1 Efficiency and grades of energy

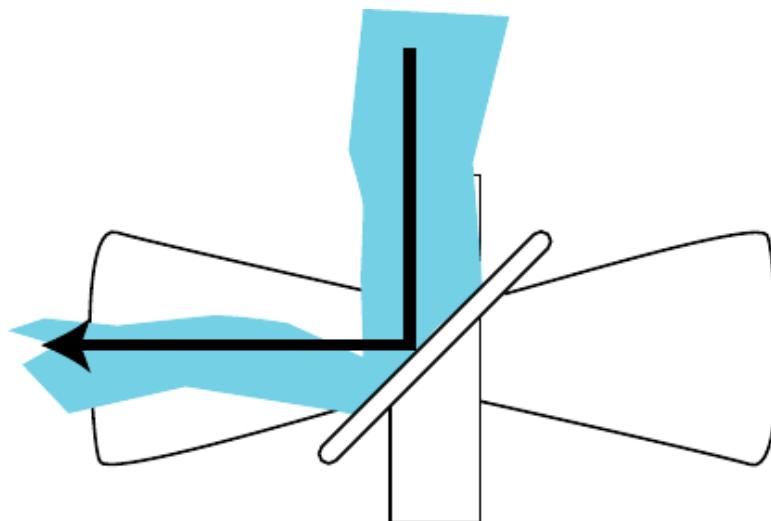
Some forms of energy are more convenient than others in certain situations. You can't run a spring-powered mechanical clock on a battery, and you can't run a battery-powered clock with mechanical energy. However, there is no fundamental physical principle that prevents you from converting 100% of the electrical energy in a battery into mechanical energy or vice-versa. More efficient motors and generators are being designed every year. In general, the laws of physics permit perfectly efficient conversion within a broad class of forms of energy.

Heat is different. Friction tends to convert other forms of energy into heat even in the best lubricated machines. When we slide a book on a table, friction brings it to a stop and converts all its kinetic energy into heat, but we never observe the opposite process, in which a book spontaneously converts heat energy into mechanical energy and starts moving! Roughly speaking, heat is different because it is disorganized. Scrambling an egg is easy. Unscrambling it is harder.

We summarize these observations by saying that heat is a lower grade of energy than other forms such as mechanical energy.

Of course it is possible to convert heat into other forms of energy such as mechanical energy, and that is what a gasoline car's engine does with the heat created by exploding the air-gas mixture. But a car engine is a tremendously inefficient device, and a great deal of the heat is simply wasted through the radiator and the exhaust. Engineers have never succeeded in creating a perfectly efficient device for converting heat energy into mechanical energy, and we now know that this is because of a deeper physical principle that is far more basic than the design of an engine.

a / 1. The temperature difference between the hot and cold parts of the air can be used to extract mechanical energy, for example with a fan blade that spins because of the rising hot air currents. 2. If the temperature of the air is first allowed to become uniform, then no mechanical energy can be extracted. The same amount of heat energy is present, but it is no longer accessible for doing mechanical work.



5.3.2 Heat engines

Heat may be more useful in some forms than in others, i.e., there are different grades of heat energy. In figure ??/1, the difference in temperature can be used to extract mechanical work with a fan blade. This principle is used in power plants, where steam is heated by burning oil or by nuclear reactions, and then allowed to expand through a turbine which has cooler steam on the other side. On a smaller scale, there is a Christmas toy, b, that consists of a small propeller spun by the hot air rising from a set of candles, very much like the setup shown in figure ??.

In figure ??/2, however, no mechanical work can be extracted because there is no difference in temperature. Although the air in ??/2 has the same total amount of energy as the air in ??/1, the heat in ??/2 is a lower grade of energy, since none of it is accessible for doing mechanical work.

In general, we define a heat engine as any device that takes heat from a reservoir of hot matter, extracts some of the heat energy to do mechanical work, and expels a lesser amount of heat into a reservoir of cold matter. The efficiency of a heat engine equals the amount of useful work extracted, W , divided by the amount of energy we had to pay for in order to heat the hot reservoir. This latter amount of heat is the same as the amount of heat the engine extracts from the high-temperature reservoir, Q_H . By conservation of energy, we have $Q_H = W + Q_L$, where Q_L is the amount of heat expelled into the low-temperature reservoir, so the efficiency of a heat engine, W/Q_H , can be rewritten as

$$\text{efficiency} = 1 - \frac{Q_L}{Q_H}. \quad [\text{efficiency of any heat engine}]$$

(As described on p. 315, we take Q_L , Q_H , and W all to be positive.)

It turns out that there is a particular type of heat engine, the Carnot engine, which, although not 100% efficient, is more efficient than any other. The grade of heat energy in a system can thus be unambiguously defined in terms of the amount of heat energy in it that cannot be extracted even by a Carnot engine.

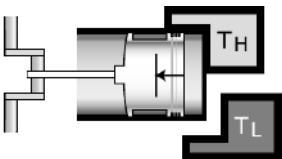
How can we build the most efficient possible engine? Let's start with an unnecessarily inefficient engine like a car engine and see how it could be improved. The radiator and exhaust expel hot gases, which is a waste of heat energy. These gases are cooler than the exploded air-gas mixture inside the cylinder, but hotter than the air that surrounds the car. We could thus improve the engine's efficiency by adding an auxiliary heat engine to it, which would operate with the first engine's exhaust as its hot reservoir and the air as its cold reservoir. In general, any heat engine that expels heat at an intermediate temperature can be made more efficient by changing it so that it expels heat only at the temperature of the cold reservoir.



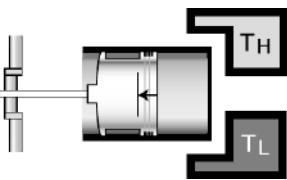
b / A heat engine. Hot air from the candles rises through the fan blades and makes the angels spin.



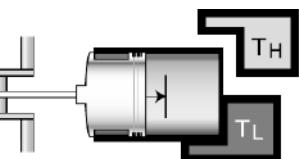
c / Sadi Carnot (1796-1832)



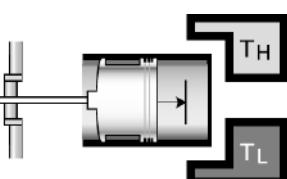
d / The beginning of the first expansion stroke, in which the working gas is kept in thermal equilibrium with the hot reservoir.



e / The beginning of the second expansion stroke, in which the working gas is thermally insulated. The working gas cools because it is doing work on the piston and thus losing energy.



f / The beginning of the first compression stroke. The working gas begins the stroke at the same temperature as the cold reservoir, and remains in thermal contact with it the whole time. The engine does negative work.



g / The beginning of the second compression stroke, in which mechanical work is absorbed, heating the working gas back up to T_H .

Similarly, any heat engine that absorbs some energy at an intermediate temperature can be made more efficient by adding an auxiliary heat engine to it which will operate between the hot reservoir and this intermediate temperature.

Based on these arguments, we define a Carnot engine as a heat engine that absorbs heat only from the hot reservoir and expels it only into the cold reservoir. Figures d-g show a realization of a Carnot engine using a piston in a cylinder filled with a monoatomic ideal gas. This gas, known as the working fluid, is separate from, but exchanges energy with, the hot and cold reservoirs. As proved on page 337, this particular Carnot engine has an efficiency given by

$$\text{efficiency} = 1 - \frac{T_L}{T_H}, \quad [\text{efficiency of a Carnot engine}]$$

where T_L is the temperature of the cold reservoir and T_H is the temperature of the hot reservoir.

Even if you do not wish to dig into the details of the proof, the basic reason for the temperature dependence is not so hard to understand. Useful mechanical work is done on strokes d and e, in which the gas expands. The motion of the piston is in the same direction as the gas's force on the piston, so positive work is done on the piston. In strokes f and g, however, the gas does negative work on the piston. We would like to avoid this negative work, but we must design the engine to perform a complete cycle. Luckily the pressures during the compression strokes are lower than the ones during the expansion strokes, so the engine doesn't undo all its work with every cycle. The ratios of the pressures are in proportion to the ratios of the temperatures, so if T_L is 20% of T_H , the engine is 80% efficient.

We have already proved that any engine that is not a Carnot engine is less than optimally efficient, and it is also true that all Carnot engines operating between a given pair of temperatures T_H and T_L have the same efficiency. (This can be proved by the methods of section 5.4.) Thus a Carnot engine is the most efficient possible heat engine.

5.3.3 Entropy

We would like to have some numerical way of measuring the grade of energy in a system. We want this quantity, called entropy, to have the following two properties:

(1) Entropy is additive. When we combine two systems and consider them as one, the entropy of the combined system equals the sum of the entropies of the two original systems. (Quantities like mass and energy also have this property.)

(2) The entropy of a system is not changed by operating a Carnot engine within it.

It turns out to be simpler and more useful to define changes in entropy than absolute entropies. Suppose as an example that a system contains some hot matter and some cold matter. It has a relatively high grade of energy because a heat engine could be used to extract mechanical work from it. But if we allow the hot and cold parts to equilibrate at some lukewarm temperature, the grade of energy has gotten worse. Thus putting heat into a hotter area is more useful than putting it into a cold area. Motivated by these considerations, we define a change in entropy as follows:

$$\Delta S = \frac{Q}{T} \quad [\text{change in entropy when adding heat } Q \text{ to matter at temperature } T; \Delta S \text{ is negative if heat is taken out}]$$

A system with a higher grade of energy has a lower entropy.

Entropy is additive.

example 10

Since changes in entropy are defined by an additive quantity (heat) divided by a non-additive one (temperature), entropy is additive.

Entropy isn't changed by a Carnot engine.

example 11

The efficiency of a heat engine is defined by

$$\text{efficiency} = 1 - Q_L/Q_H,$$

and the efficiency of a Carnot engine is

$$\text{efficiency} = 1 - T_L/T_H,$$

so for a Carnot engine we have $Q_L/Q_H = T_L/T_H$, which can be rewritten as $Q_L/T_L = Q_H/T_H$. The entropy lost by the hot reservoir is therefore the same as the entropy gained by the cold one.

Entropy increases in heat conduction.

example 12

When a hot object gives up energy to a cold one, conservation of energy tells us that the amount of heat lost by the hot object is the same as the amount of heat gained by the cold one. The change in entropy is $-Q/T_H + Q/T_L$, which is positive because $T_L < T_H$.

Entropy is increased by a non-Carnot engine.

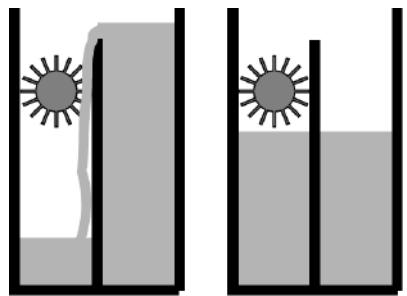
example 13

The efficiency of a non-Carnot engine is less than $1 - T_L/T_H$, so $Q_L/Q_H > T_L/T_H$ and $Q_L/T_L > Q_H/T_H$. This means that the entropy increase in the cold reservoir is greater than the entropy decrease in the hot reservoir.

A book sliding to a stop

example 14

A book slides across a table and comes to a stop. Once it stops, all its kinetic energy has been transformed into heat. As the book and table heat up, their entropies both increase, so the total entropy increases as well.



h / Entropy can be understood using the metaphor of a water wheel. Letting the water levels equalize is like letting the entropy maximize. Taking water from the high side and putting it into the low side increases the entropy. Water levels in this metaphor correspond to temperatures in the actual definition of entropy.

All of these examples involved closed systems, and in all of them the total entropy either increased or stayed the same. It never decreased. Here are two examples of schemes for decreasing the entropy of a closed system, with explanations of why they don't work.

Using a refrigerator to decrease entropy? example 15

▷ A refrigerator takes heat from a cold area and dumps it into a hot area. (1) Does this lead to a net decrease in the entropy of a closed system? (2) Could you make a Carnot engine more efficient by running a refrigerator to cool its low-temperature reservoir and eject heat into its high-temperature reservoir?

▷ (1) No. The heat that comes off of the radiator coils is a great deal more than the heat the fridge removes from inside; the difference is what it costs to run your fridge. The heat radiated from the coils is so much more than the heat removed from the inside that the increase in the entropy of the air in the room is greater than the decrease of the entropy inside the fridge. The most efficient refrigerator is actually a Carnot engine running in reverse, which leads to neither an increase nor a decrease in entropy.

(2) No. The most efficient refrigerator is a reversed Carnot engine. You will not achieve anything by running one Carnot engine in reverse and another forward. They will just cancel each other out.

Maxwell's demon example 16

▷ Maxwell imagined a pair of rooms, their air being initially in thermal equilibrium, having a partition across the middle with a tiny door. A minuscule demon is posted at the door with a little ping-pong paddle, and his duty is to try to build up faster-moving air molecules in room B and slower moving ones in room A. For instance, when a fast molecule is headed through the door, going from A to B, he lets it by, but when a slower than average molecule tries the same thing, he hits it back into room A. Would this decrease the total entropy of the pair of rooms?

▷ No. The demon needs to eat, and we can think of his body as a little heat engine, and his metabolism is less efficient than a Carnot engine, so he ends up increasing the entropy rather than decreasing it.

Observations such as these lead to the following hypothesis, known as the second law of thermodynamics:

The entropy of a closed system always increases, or at best stays the same: $\Delta S \geq 0$.

At present our arguments to support this statement may seem less than convincing, since they have so much to do with obscure facts about heat engines. In the following section we will find a more satisfying and fundamental explanation for the continual increase in

entropy. To emphasize the fundamental and universal nature of the second law, here are a few exotic examples.

Entropy and evolution

example 17

A favorite argument of many creationists who don't believe in evolution is that evolution would violate the second law of thermodynamics: the death and decay of a living thing releases heat (as when a compost heap gets hot) and lessens the amount of energy available for doing useful work, while the reverse process, the emergence of life from nonliving matter, would require a decrease in entropy. Their argument is faulty, since the second law only applies to closed systems, and the earth is not a closed system. The earth is continuously receiving energy from the sun.

The heat death of the universe

example 18

Living things have low entropy: to demonstrate this fact, observe how a compost pile releases heat, which then equilibrates with the cooler environment. We never observe dead things to leap back to life after sucking some heat energy out of their environments! The only reason life was able to evolve on earth was that the earth was not a closed system: it got energy from the sun, which presumably gained more entropy than the earth lost.

Victorian philosophers spent a lot of time worrying about the heat death of the universe: eventually the universe would have to become a high-entropy, lukewarm soup, with no life or organized motion of any kind. Fortunately (?), we now know a great many other things that will make the universe inhospitable to life long before its entropy is maximized. Life on earth, for instance, will end when the sun evolves into a hotter state and boils away our oceans.

Hawking radiation

example 19

Any process that could destroy heat (or convert it into nothing but mechanical work) would lead to a reduction in entropy. Black holes are supermassive stars whose gravity is so strong that nothing, not even light, can escape from them once it gets within a boundary known as the event horizon. Black holes are commonly observed to suck hot gas into them. Does this lead to a reduction in the entropy of the universe? Of course one could argue that the entropy is still there inside the black hole, but being able to "hide" entropy there amounts to the same thing as being able to destroy entropy.

The physicist Steven Hawking was bothered by this, and finally realized that although the actual stuff that enters a black hole is lost forever, the black hole will gradually lose energy in the form of light emitted from just outside the event horizon. This light ends up reintroducing the original entropy back into the universe.

Discussion Questions

A In this discussion question, you'll think about a car engine in terms of thermodynamics. Note that an internal combustion engine doesn't fit very well into the theoretical straightjacket of a heat engine. For instance, a heat engine has a high-temperature heat reservoir at a single well-defined temperature, T_H . In a typical car engine, however, there are several very different temperatures you could imagine using for T_H : the temperature of the engine block ($\sim 100^\circ\text{C}$), the walls of the cylinder ($\sim 250^\circ\text{C}$), or the temperature of the exploding air-gas mixture ($\sim 1000^\circ\text{C}$, with significant changes over a four-stroke cycle). Let's use $T_H \sim 1000^\circ\text{C}$.

Burning gas supplies heat energy Q_H to your car's engine. The engine does mechanical work W , but also expels heat Q_L into the environment through the radiator and the exhaust. Conservation of energy gives

$$Q_H = Q_L + W,$$

and the relative proportions of Q_L and W are usually about 90% to 10%. (Actually it depends quite a bit on the type of car, the driving conditions, etc.) Here, Q_H , Q_L , and W are all positive according to the sign convention defined on p. 315.

(1) A gallon of gas releases about 140 MJ of heat Q_H when burned. Estimate the change in entropy of the universe due to running a typical car engine and burning one gallon of gas. Note that you'll have to introduce appropriate plus and minus signs, as defined in the relation $\Delta S = Q/T$, in which heat input raises an object's entropy and heat output lowers it. (You'll have to estimate how hot the environment is. For the sake of argument, assume that the work done by the engine, W , remains in the form of mechanical energy, although in reality it probably ends up being changed into heat when you step on the brakes.) Is your result consistent with the second law of thermodynamics?

(2) Q_L is obviously undesirable: you pay for it, but all it does is heat the neighborhood. Suppose that engineers do a really good job of getting rid of the effects that create Q_L , such as friction. Could Q_L ever be reduced to zero, at least theoretically? What would happen if you redid the calculation in #1, but assumed $Q_L = 0$?

B When we run the Carnot engine in figures d-g, there are four parts of the universe that undergo changes in their physical states: the hot reservoir, the cold reservoir, the working gas, and the outside world to which the shaft is connected in order to do physical work. Over one full cycle, discuss which of these parts gain entropy, which ones lose entropy, and which ones keep the same entropy. During which of the four strokes do these changes occur?

5.4 Entropy as a microscopic quantity

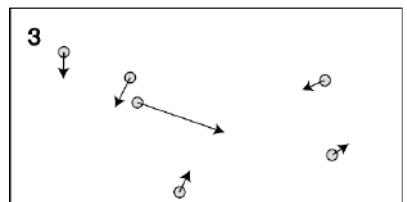
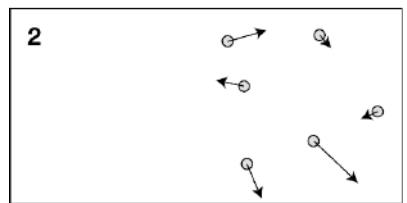
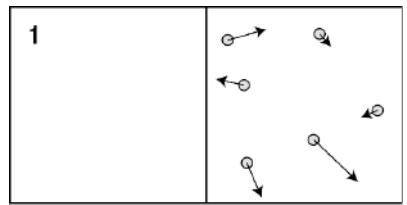
5.4.1 A microscopic view of entropy

To understand why the second law of thermodynamics is always true, we need to see what entropy really means at the microscopic level. An example that is easy to visualize is the free expansion of a monoatomic gas. Figure a/1 shows a box in which all the atoms of the gas are confined on one side. We very quickly remove the

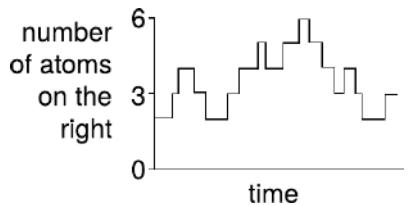
barrier between the two sides, $a/2$, and some time later, the system has reached an equilibrium, $a/3$. Each snapshot shows both the positions and the momenta of the atoms, which is enough information to allow us in theory to extrapolate the behavior of the system into the future, or the past. However, with a realistic number of atoms, rather than just six, this would be beyond the computational power of any computer.²

But suppose we show figure a/2 to a friend without any further information, and ask her what she can say about the system's behavior in the future. She doesn't know how the system was prepared. Perhaps, she thinks, it was just a strange coincidence that all the atoms happened to be in the right half of the box at this particular moment. In any case, she knows that this unusual situation won't last for long. She can predict that after the passage of any significant amount of time, a surprise inspection is likely to show roughly half the atoms on each side. The same is true if you ask her to say what happened in the past. She doesn't know about the barrier, so as far as she's concerned, extrapolation into the past is exactly the same kind of problem as extrapolation into the future. We just have to imagine reversing all the momentum vectors, and then all our reasoning works equally well for backwards extrapolation. She would conclude, then, that the gas in the box underwent an unusual fluctuation, b, and she knows that the fluctuation is very unlikely to exist very far into the future, or to have existed very far into the past.

What does this have to do with entropy? Well, state a/3 has a greater entropy than state a/2. It would be easy to extract mechanical work from a/2, for instance by letting the gas expand while pressing on a piston rather than simply releasing it suddenly into the void. There is no way to extract mechanical work from state a/3. Roughly speaking, our microscopic description of entropy relates to the *number of possible states*. There are a lot more states like a/3 than there are states like a/2. Over long enough periods of time — long enough for equilibration to occur — the system gets mixed up, and is about equally likely to be in any of its possible states, regardless of what state it was initially in. We define some number that describes an interesting property of the whole system, say the number of atoms in the right half of the box, R . A high-entropy value of R is one like $R = 3$, which allows many possible states. We are far more likely to encounter $R = 3$ than a low-entropy value like $R = 0$ or $R = 6$.



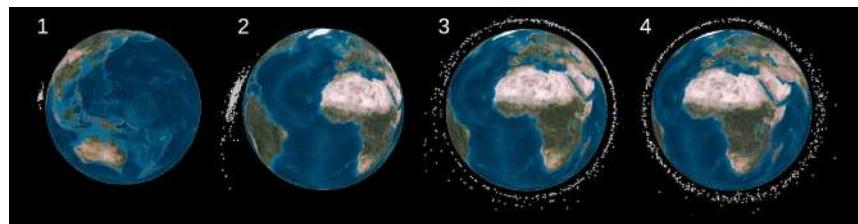
a / A gas expands freely, doubling its volume.



b / An unusual fluctuation in the distribution of the atoms between the two sides of the box. There has been no external manipulation as in figure a/1.

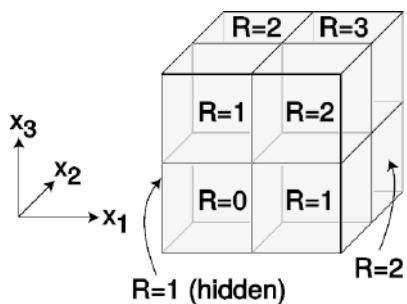
²Even with smaller numbers of atoms, there is a problem with this kind of brute-force computation, because the tiniest measurement errors in the initial state would end up having large effects later on.

c / Earth orbit is becoming cluttered with space junk, and the pieces can be thought of as the “molecules” comprising an exotic kind of gas. These images show the evolution of a cloud of debris arising from a 2007 Chinese test of an anti-satellite rocket. Panels 1-4 show the cloud five minutes, one hour, one day, and one month after the impact. The entropy seems to have maximized by panel 4.



	$R=1$	$R=2$
x_2		
	$R=0$	$R=1$
x_1		

d / The phase space for two atoms in a box.



e / The phase space for three atoms in a box.

5.4.2 Phase space

There is a problem with making this description of entropy into a mathematical definition. The problem is that it refers to the number of possible states, but that number is theoretically infinite. To get around the problem, we coarsen our description of the system. For the atoms in figure a, we don’t really care exactly where each atom is. We only care whether it is in the right side or the left side. If a particular atom’s left-right position is described by a coordinate x , then the set of all possible values of x is a line segment along the x axis, containing an infinite number of points. We break this line segment down into two halves, each of width Δx , and we consider two different values of x to be variations on the same state if they both lie in the same half. For our present purposes, we can also ignore completely the y and z coordinates, and all three momentum components, p_x , p_y , and p_z .

Now let’s do a real calculation. Suppose there are only two atoms in the box, with coordinates x_1 and x_2 . We can give all the relevant information about the state of the system by specifying one of the cells in the grid shown in figure d. This grid is known as the *phase space* of the system.³ The lower right cell, for instance, describes a state in which atom number 1 is in the right side of the box and atom number 2 in the left. Since there are two possible states with $R = 1$ and only one state with $R = 2$, we are twice as likely to observe $R = 1$, and $R = 1$ has higher entropy than $R = 2$.

Figure e shows a corresponding calculation for three atoms, which makes the phase space three-dimensional. Here, the $R = 1$ and 2 states are three times more likely than $R = 0$ and 3. Four atoms would require a four-dimensional phase space, which exceeds our ability to visualize. Although our present example doesn’t require it, a phase space can describe momentum as well as position, as shown in figure f. In general, a phase space for a monoatomic gas has six dimensions per atom (one for each coordinate and one for each momentum component).

³The term is a little obscure. Basically the idea is the same as in “my toddler is going through a phase where he always says no.” The “phase” is a stage in the evolution of the system, a snapshot of its state at a moment in time. The usage is also related to the concept of Lissajous figures, in which a particular point on the trajectory is defined by the phases of the oscillations along the x and y axes.

5.4.3 Microscopic definitions of entropy and temperature

Two more issues need to be resolved in order to make a microscopic definition of entropy.

First, if we defined entropy as the number of possible states, it would be a multiplicative quantity, not an additive one: if an ice cube in a glass of water has M_1 states available to it, and the number of states available to the water is M_2 , then the number of possible states of the whole system is the product $M_1 M_2$. To get around this problem, we take the natural logarithm of the number of states, which makes the entropy additive because of the property of the logarithm $\ln(M_1 M_2) = \ln M_1 + \ln M_2$.

The second issue is a more trivial one. The concept of entropy was originally invented as a purely macroscopic quantity, and the macroscopic definition $\Delta S = Q/T$, which has units of J/K, has a different calibration than would result from defining $S = \ln M$. The calibration constant we need turns out to be simply the Boltzmann constant, k .

Microscopic definition of entropy: The entropy of a system is $S = k \ln M$, where M is the number of available states.⁴

This also leads to a more fundamental definition of temperature. Two systems are in thermal equilibrium when they have maximized their combined entropy through the exchange of energy. Here the energy possessed by one part of the system, E_1 or E_2 , plays the same role as the variable R in the examples of free expansion above. A maximum of a function occurs when the derivative is zero, so the maximum entropy occurs when

$$\frac{d(S_1 + S_2)}{dE_1} = 0.$$

We assume the systems are only able to exchange heat energy with each other, $dE_1 = -dE_2$, so

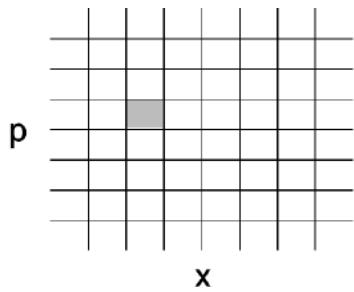
$$\frac{dS_1}{dE_1} = \frac{dS_2}{dE_2},$$

and since the energy is being exchanged in the form of heat we can make the equations look more familiar if we write dQ for an amount of heat to be transferred into either system:

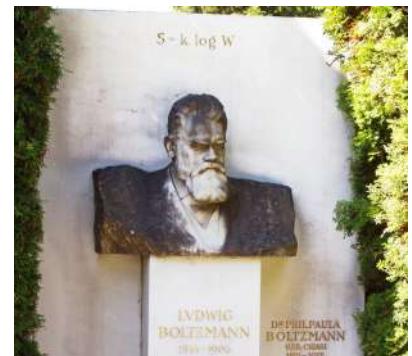
$$\frac{dS_1}{dQ_1} = \frac{dS_2}{dQ_2}.$$

In terms of our previous definition of entropy, this is equivalent to $1/T_1 = 1/T_2$, which makes perfect sense since the systems are in thermal equilibrium. According to our new approach, entropy has

⁴This is the same relation as the one on Boltzmann's tomb, just in a slightly different notation.



f / A phase space for a single atom in one dimension, taking momentum into account.



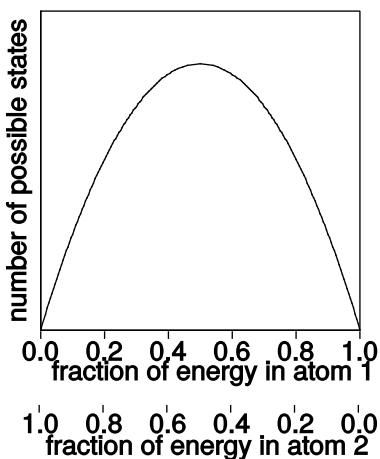
g / Ludwig Boltzmann's tomb, inscribed with his equation for entropy.

already been defined in a fundamental manner, so we can take this as a definition of temperature:

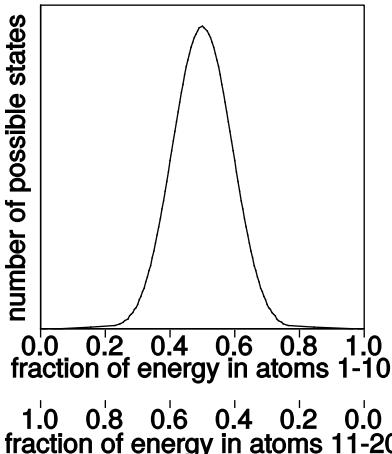
$$\frac{1}{T} = \frac{dS}{dQ},$$

where dS represents the increase in the system's entropy from adding heat dQ to it.

Examples with small numbers of atoms



h / A two-atom system has the highest number of available states when the energy is equally divided. Equal energy division is therefore the most likely possibility at any given moment in time.



i / When two systems of 10 atoms each interact, the graph of the number of possible states is narrower than with only one atom in each system.

Let's see how this applies to an ideal, monoatomic gas with a small number of atoms. To start with, consider the phase space available to one atom. Since we assume the atoms in an ideal gas are noninteracting, their positions relative to each other are really irrelevant. We can therefore enumerate the number of states available to each atom just by considering the number of momentum vectors it can have, without considering its possible locations. The relationship between momentum and kinetic energy is $E = (p_x^2 + p_y^2 + p_z^2)/2m$, so if for a fixed value of its energy, we arrange all of an atom's possible momentum vectors with their tails at the origin, their tips all lie on the surface of a sphere in phase space with radius $|\mathbf{p}| = \sqrt{2mE}$. The number of possible states for that atom is proportional to the sphere's surface area, which in turn is proportional to the square of the sphere's radius, $|\mathbf{p}|^2 = 2mE$.

Now consider two atoms. For any given way of sharing the energy between the atoms, $E = E_1 + E_2$, the number of possible combinations of states is proportional to $E_1 E_2$. The result is shown in figure *h*. The greatest number of combinations occurs when we divide the energy equally, so an equal division gives maximum entropy.

By increasing the number of atoms, we get a graph whose peak is narrower, i. With more than one atom in each system, the total energy is $E = (p_{x,1}^2 + p_{y,1}^2 + p_{z,1}^2 + p_{x,2}^2 + p_{y,2}^2 + p_{z,2}^2 + \dots)/2m$. With n atoms, a total of $3n$ momentum coordinates are needed in order to specify their state, and such a set of numbers is like a single point in a $3n$ -dimensional space (which is impossible to visualize). For a given total energy E , the possible states are like the surface of a $3n$ -dimensional sphere, with a surface area proportional to p^{3n-1} , or $E^{(3n-1)/2}$. The graph in figure *i*, for example, was calculated according to the formula $E_1^{29/2} E_2^{29/2} = E_1^{29/2} (E - E_1)^{29/2}$.

Since graph *i* is narrower than graph *h*, the fluctuations in energy sharing are smaller. If we inspect the system at a random moment in time, the energy sharing is very unlikely to be more lopsided than a 40-60 split. Now suppose that, instead of 10 atoms interacting with 10 atoms, we had a 10^{23} atoms interacting with 10^{23} atoms. The graph would be extremely narrow, and it would be a statistical certainty that the energy sharing would be nearly perfectly equal. This is why we never observe a cold glass of water to change itself

into an ice cube sitting in some warm water!

By the way, note that although we've redefined temperature, these examples show that things are coming out consistent with the old definition, since we saw that the old definition of temperature could be described in terms of the average energy per atom, and here we're finding that equilibration results in each subset of the atoms having an equal share of the energy.

Entropy of a monoatomic ideal gas

Let's calculate the entropy of a monoatomic ideal gas of n atoms. This is an important example because it allows us to show that our present microscopic treatment of thermodynamics is consistent with our previous macroscopic approach, in which temperature was defined in terms of an ideal gas thermometer.

The number of possible locations for each atom is $V/\Delta x^3$, where Δx is the size of the space cells in phase space. The number of possible combinations of locations for the atoms is therefore $(V/\Delta x^3)^n$.

The possible momenta cover the surface of a $3n$ -dimensional sphere, whose radius is $\sqrt{2mE}$, and whose surface area is therefore proportional to $E^{(3n-1)/2}$. In terms of phase-space cells, this area corresponds to $E^{(3n-1)/2}/\Delta p^{3n}$ possible combinations of momenta, multiplied by some constant of proportionality which depends on m , the atomic mass, and n , the number of atoms. To avoid having to calculate this constant of proportionality, we limit ourselves to calculating the part of the entropy that does not depend on n , so the resulting formula will not be useful for comparing entropies of ideal gas samples with different numbers of atoms.

The final result for the number of available states is

$$M = \left(\frac{V}{\Delta x^3} \right)^n \frac{E^{(3n-1)/2}}{\Delta p^{3n}}, \quad [\text{function of } n]$$

so the entropy is

$$S = nk \ln V + \frac{3}{2}nk \ln E + (\text{function of } \Delta x, \Delta p, \text{ and } n),$$

where the distinction between n and $n - 1$ has been ignored. Using $PV = nkT$ and $E = (3/2)nkT$, we can also rewrite this as

$$S = \frac{5}{2}nk \ln T - nk \ln P + \dots, \quad [\text{entropy of a monoatomic ideal gas}]$$

where “ \dots ” indicates terms that may depend on Δx , Δp , m , and n , but that have no effect on comparisons of gas samples with the same number of atoms.

self-check C

Why does it make sense that the temperature term has a positive sign in the above example, while the pressure term is negative? Why does

it make sense that the whole thing is proportional to n ? ▷ Answer, p. 1061

To show consistency with the macroscopic approach to thermodynamics, we need to show that these results are consistent with the behavior of an ideal-gas thermometer. Using the new definition $1/T = dS/dQ$, we have $1/T = dS/dE$, since transferring an amount of heat dQ into the gas increases its energy by a corresponding amount. Evaluating the derivative, we find $1/T = (3/2)nk/E$, or $E = (3/2)nkT$, which is the correct relation for a monoatomic ideal gas.

A mixture of molecules

example 20

▷ Suppose we have a mixture of two different monoatomic gases, say helium and argon. How would we find the entropy of such a mixture (say, in terms of V and E)? How would the energy be shared between the two types of molecules, i.e., would a more massive argon atom have more energy on the average than a less massive helium atom, the same, or less?

▷ Since entropy is additive, we simply need to add the entropies of the two types of atom. However, the expression derived above for the entropy omitted the dependence on the mass m of the atom, which is different for the two constituents of the gas, so we need to go back and figure out how to put that m -dependence back in. The only place where we threw away m 's was when we identified the radius of the sphere in momentum space with $\sqrt{2mE}$, but then threw away the constant factor of m . In other words, the final result can be generalized merely by replacing E everywhere with the product mE . Since the log of a product is the sum of the logs, the dependence of the final result on m and E can be broken apart into two different terms, and we find

$$S = nk \ln V + \frac{3}{2}nk \ln m + \frac{3}{2}nk \ln E + \dots$$

The total entropy of the mixture can then be written as

$$\begin{aligned} S = n_1 k \ln V + n_2 k \ln V + \frac{3}{2}n_1 k \ln m_1 + \frac{3}{2}n_2 k \ln m_2 \\ + \frac{3}{2}n_1 k \ln E_1 + \frac{3}{2}n_2 k \ln E_2 + \dots \end{aligned}$$

Now what about the energy sharing? If the total energy is $E = E_1 + E_2$, then the most overwhelmingly probable sharing of energy will be the one that maximizes the entropy. Notice that the dependence of the entropy on the masses m_1 and m_2 occurs in terms that are entirely separate from the energy terms. If we want to maximize S with respect to E_1 (with $E_2 = E - E_1$ by conservation of energy), then we differentiate S with respect to E_1 and set it equal to zero. The terms that contain the masses don't have any

dependence on E_1 , so their derivatives are zero, and we find that the molecular masses can have no effect on the energy sharing. Setting the derivative equal to zero, we have

$$\begin{aligned} 0 &= \frac{\partial}{\partial E_1} \left(n_1 k \ln V + n_2 k \ln V + \frac{3}{2} n_1 k \ln m_1 + \frac{3}{2} n_2 k \ln m_2 \right. \\ &\quad \left. + \frac{3}{2} n_1 k \ln E_1 + \frac{3}{2} n_2 k \ln(E - E_1) + \dots \right) \\ &= \frac{3}{2} k \left(\frac{n_1}{E_1} - \frac{n_2}{E - E_1} \right) \\ 0 &= \frac{n_1}{E_1} - \frac{n_2}{E - E_1} \\ \frac{n_1}{E_1} &= \frac{n_2}{E_2}. \end{aligned}$$

In other words, each gas gets a share of the energy in proportion to the number of its atoms, and therefore every atom gets, on average, the same amount of energy, regardless of its mass. The result for the average energy per atom is exactly the same as for an unmixed gas, $\bar{K} = (3/2)kT$.

5.4.4 Equipartition

Betsy Salazar of Redwood Cove, California, has 37 pet raccoons, which is theoretically illegal. She admits that she has trouble telling them apart, but she tries to give them all plenty of care and affection (which they reciprocate). There are only so many hours in a day, so there is a fixed total amount of love. The raccoons share this love unequally on any given day, but *on the average* they all get the same amount. This kind of equal-sharing-on-the-average-out-of-some-total-amount is more concisely described using the term *equipartition*, meaning equal partitioning, or equal sharing. If Betsy did keep track of how much love she lavished on each animal, using some numerical scale, we would have 37 numbers to keep track of. We say that the love is partitioned among 37 *degrees of freedom*.

For a monoatomic ideal gas, the analysis in section 5.2.2, p. 317, leads to the following simple and useful fact about how kinetic energy is shared among all the atoms' x , y , and z degrees of freedom:

Equipartition theorem: restricted form

For a monoatomic ideal gas containing n atoms, each of the $3n$ degrees of freedom contains, on the average, a kinetic energy $\frac{1}{2}kT$.

As applications of this fact, we can easily find the amount of heat needed to raise the temperature of the gas by one unit (its specific heat), or estimate the typical thermal velocities of the atoms given the temperature. Equipartition gives us a more rigorous and quantitative statement of the idea that temperature is a measure of how concentrated the heat is, or of how much energy there is per particle.

We would now like to generalize this theorem. Example 20, p. 332, tells us that it doesn't matter whether the gas is a mixture of two different types of atoms — Betsy doesn't give a raccoon more love just because it's big and fat. Equal energy sharing might seem obvious by symmetry if all the atoms are identical, but we see that it still holds when they are not identical. Symmetry was not a necessary assumption.

To generalize even further, let's look at what the necessary assumptions really were in example 20. For simplicity, suppose we have only one argon atom, named Alice, and one helium atom, named Harry. Their total kinetic energy is $E = p_x^2/2m + p_y^2/2m + p_z^2/2m + p'^2_x/2m' + p'^2_y/2m' + p'^2_z/2m'$, where the primes indicate Harry. The system consisting of Alice and Harry has six degrees of freedom (the six momenta), and the six terms in the energy all look alike. The only difference among them is that the constant factors attached to the squares of the momenta have different values, but we've just proved that those differences don't matter. In other words, if we have any system at all whose energy is of the form $E = (\dots)p_1^2 + (\dots)p_2^2 + \dots$, with any number of terms, then each term holds, on average, the same amount of energy, $\frac{1}{2}kT$.

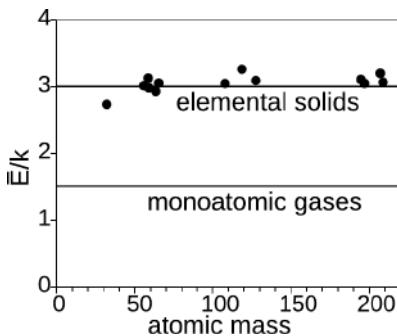
It doesn't even matter whether the things being squared are momenta: if you look back over the logical steps that went into the argument, you'll see that none of them depended on that. In a solid, for example, the atoms aren't free to wander around, but they can vibrate from side to side. If an atom moves away from its equilibrium position at $x = 0$ to some other value of x , then its electrical energy is $(1/2)\kappa x^2$, where κ is the spring constant (written as the Greek letter kappa to distinguish it from the Boltzmann constant k). We can conclude that each atom in the solid, on average, has $\frac{1}{2}kT$ of energy in the electrical energy due to its x displacement along the x axis, and equal amounts for y and z . Thus for a solid, we expect there to be a total of *six* energies per atom (three kinetic energies and three interaction energies), each of which carries an average energy $\frac{1}{2}kT$, for a total of $3kT$. In other words, a solid should have twice the heat capacity of a monoatomic gas. This was discovered empirically by Dulong and Petit in 1819, as shown in figure j.

Equipartition theorem: general form

For a system whose energy can be written as the sum of the squares of n variables, the average value of each term in the energy is $\frac{1}{2}kT$.

An unexpected glimpse of the microcosm

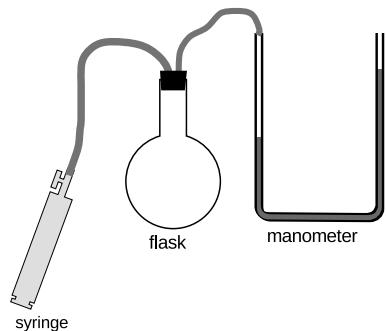
These ideas about equipartition now lead us to some surprising insights into how the microscopic world manifests itself on the human scale. You may have the feeling at this point that of course Boltzmann was right about the literal existence of atoms, but only



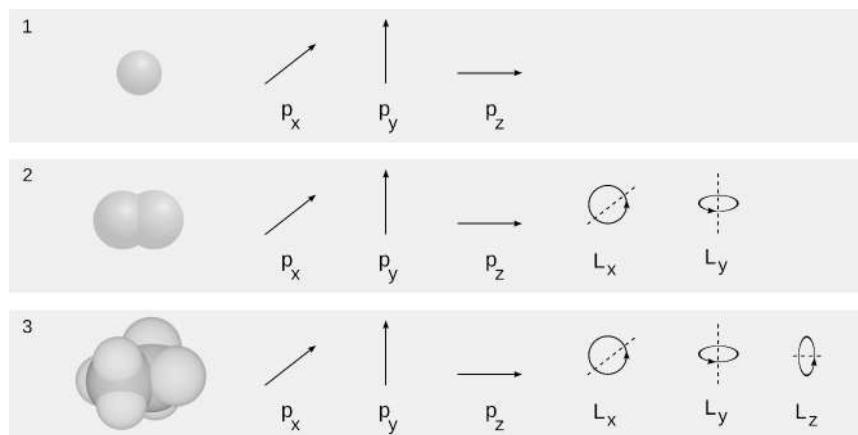
j / Heat capacities of solids cluster around $3kT$ per atom. The elemental solids plotted are the ones originally used by Dulong and Petit to infer empirically that the heat capacity of solids per atom was constant. (Modern data.)

very sophisticated experiments could vindicate him definitively. After all, the microscopic and macroscopic definitions of entropy are equivalent, so it might seem as though there was no real advantage to the microscopic approach. Surprisingly, very simple experiments are capable of revealing a picture of the microscopic world, and there is no possible macroscopic explanation for their results.

In 1819, before Boltzmann was born, Clément and Desormes did an experiment like the one shown in figure k. The gas in the flask is pressurized using the syringe. This heats it slightly, so it is then allowed to cool back down to room temperature. Its pressure is measured using the manometer. The stopper on the flask is popped and then immediately reinserted. Its pressure is now equalized with that in the room, and the gas's expansion has cooled it a little, because it did mechanical work on its way out of the flask, causing it to lose some of its internal energy E . The expansion is carried out quickly enough so that there is not enough time for any significant amount of heat to flow in through the walls of the flask before the stopper is reinserted. The gas is now allowed to come back up to room temperature (which takes a much longer time), and as a result regains a fraction b of its original overpressure. During this constant-volume reheating, we have $PV = nkT$, so the amount of pressure regained is a direct indication of how much the gas cooled down when it lost an amount of energy ΔE .



k / An experiment for determining the shapes of molecules.



I / The differing shapes of a helium atom (1), a nitrogen molecule (2), and a difluoroethane molecule (3) have surprising macroscopic effects.

If the gas is monoatomic, then we know what to expect for this relationship between energy and temperature: $\Delta E = (3/2)nk\Delta T$, where the factor of 3 came ultimately from the fact that the gas was in a three-dimensional space, 1/1. Moving in this space, each molecule can have momentum in the x , y , and z directions. It has three degrees of freedom. What if the gas is not monoatomic? Air, for example, is made of diatomic molecules, 1/2. There is a subtle difference between the two cases. An individual atom of a monoatomic gas is a perfect sphere, so it is exactly the same no

matter how it is oriented. Because of this perfect symmetry, there is thus no way to tell whether it is spinning or not, and in fact we find that it can't rotate. The diatomic gas, on the other hand, can rotate end over end about the x or y axis, but cannot rotate about the z axis, which is its axis of symmetry. It has a total of five degrees of freedom. A polyatomic molecule with a more complicated, asymmetric shape, $1/3$, can rotate about all three axis, so it has a total of six degrees of freedom.

Because a polyatomic molecule has more degrees of freedom than a monoatomic one, it has more possible states for a given amount of energy. That is, its entropy is higher for the same energy. From the definition of temperature, $1/T = dS/dE$, we conclude that it has a lower temperature for the same energy. In other words, it is more difficult to heat n molecules of difluoroethane than it is to heat n atoms of helium. When the Clément-Desormes experiment is carried out, the result b therefore depends on the shape of the molecule! Who would have dreamed that such simple observations, correctly interpreted, could give us this kind of glimpse of the microcosm?

Lets go ahead and calculate how this works. Suppose a gas is allowed to expand without being able to exchange heat with the rest of the universe. The loss of thermal energy from the gas equals the work it does as it expands, and using the result of homework problem 2 on page 347, the work done in an infinitesimal expansion equals $P dV$, so

$$dE + P dV = 0.$$

(If the gas had not been insulated, then there would have been a third term for the heat gained or lost by heat conduction.)

From section 5.2 we have $E = (3/2)PV$ for a monoatomic ideal gas. More generally, the equipartition theorem tells us that the 3 simply needs to be replaced with the number of degrees of freedom α , so $dE = (\alpha/2)P dV + (\alpha/2)V dP$, and the equation above becomes

$$0 = \frac{\alpha+2}{2}P dV + \frac{\alpha}{2}V dP.$$

Rearranging, we have

$$(\alpha+2)\frac{dV}{V} = -\alpha\frac{dP}{P}.$$

Integrating both sides gives

$$(\alpha+2)\ln V = -\alpha\ln P + \text{constant},$$

and taking exponentials on both sides yields

$$V^{\alpha+2} \propto P^{-\alpha}.$$

We now wish to reexpress this in terms of pressure and temperature. Eliminating $V \propto (T/P)$ gives

$$T \propto P^b,$$

where $b = 2/(\alpha+2)$ is equal to $2/5$, $2/7$, or $1/4$, respectively, for a monoatomic, diatomic, or polyatomic gas.

Efficiency of the Carnot engine

example 21

As an application, we now prove the result claimed earlier for the efficiency of a Carnot engine. First consider the work done during the constant-temperature strokes. Integrating the equation $dW = PdV$, we have $W = \int PdV$. Since the thermal energy of an ideal gas depends only on its temperature, there is no change in the thermal energy of the gas during this constant-temperature process. Conservation of energy therefore tells us that work done by the gas must be exactly balanced by the amount of heat transferred in from the reservoir.

$$\begin{aligned} Q &= W \\ &= \int PdV \end{aligned}$$

For our proof of the efficiency of the Carnot engine, we need only the ratio of Q_H to Q_L , so we neglect constants of proportionality, and simply substitute $P \propto T/V$, giving

$$Q \propto \int \frac{T}{V} dV \propto T \ln \frac{V_2}{V_1} \propto T \ln \frac{P_1}{P_2}.$$

The efficiency of a heat engine is

$$\text{efficiency} = 1 - \frac{Q_L}{Q_H}.$$

Making use of the result from the previous proof for a Carnot engine with a monoatomic ideal gas as its working gas, we have

$$\text{efficiency} = 1 - \frac{T_L \ln(P_4/P_3)}{T_H \ln(P_1/P_2)},$$

where the subscripts 1, 2, 3, and 4 refer to figures d–g on page 322. We have shown above that the temperature is proportional to P^b on the insulated strokes 2–3 and 4–1, the pressures must be related by $P_2/P_3 = P_1/P_4$, which can be rearranged as $P_4/P_3 = P_1/P_2$, and we therefore have

$$\text{efficiency} = 1 - \frac{T_L}{T_H}.$$

5.4.5 The arrow of time, or “this way to the Big Bang”

Now that we have a microscopic understanding of entropy, what does that tell us about the second law of thermodynamics? The second law defines a forward direction to time, “time’s arrow.” The microscopic treatment of entropy, however, seems to have mysteriously sidestepped that whole issue. A graph like figure b on page 327, showing a fluctuation away from equilibrium, would look just

as natural if we flipped it over to reverse the direction of time. After all, the basic laws of physics are conservation laws, which don't distinguish between past and future. Our present picture of entropy suggests that we restate the second law of thermodynamics as follows: low-entropy states are short-lived. An ice cube can't exist forever in warm water. We no longer have to distinguish past from future.

But how do we reconcile this with our strong psychological sense of the direction of time, including our ability to remember the past but not the future? Why do we observe ice cubes melting in water, but not the time-reversed version of the same process?

The answer is that there is no past-future asymmetry in the laws of physics, but there is a past-future asymmetry in the universe. The universe started out with the Big Bang. (Some of the evidence for the Big Bang theory is given on page 370.) The early universe had a very low entropy, and low-entropy states are short-lived. What does "short-lived" mean here, however? Hot coffee left in a paper cup will equilibrate with the air within ten minutes or so. Hot coffee in a thermos bottle maintains its low-entropy state for much longer, because the coffee is insulated by a vacuum between the inner and outer walls of the thermos. The universe has been mostly vacuum for a long time, so it's well insulated. Also, it takes billions of years for a low-entropy normal star like our sun to evolve into the high-entropy cinder known as a white dwarf.

The universe, then, is still in the process of equilibrating, and all the ways we have of telling the past from the future are really just ways of determining which direction in time points toward the Big Bang, i.e., which direction points to lower entropy. The psychological arrow of time, for instance, is ultimately based on the thermodynamic arrow. In some general sense, your brain is like a computer, and computation has thermodynamic effects. In even the most efficient possible computer, for example, erasing one bit of memory decreases its entropy from $k \ln 2$ (two possible states) to $k \ln 1$ (one state), for a drop of about 10^{-23} J/K . One way of determining the direction of the psychological arrow of time is that forward in psychological time is the direction in which, billions of years from now, all consciousness will have ceased; if consciousness was to exist forever in the universe, then there would have to be a never-ending decrease in the universe's entropy. This can't happen, because low-entropy states are short-lived.

Relating the direction of the thermodynamic arrow of time to the existence of the Big Bang is a satisfying way to avoid the paradox of how the second law can come from basic laws of physics that don't distinguish past from future. There is a remaining mystery, however: why did our universe have a Big Bang that was low in entropy? It could just as easily have been a maximum-entropy state, and in fact

the number of possible high-entropy Big Bangs is vastly greater than the number of possible low-entropy ones. The question, however, is probably not one that can be answered using the methods of science. All we can say is that if the universe had started with a maximum-entropy Big Bang, then we wouldn't be here to wonder about it. A longer, less mathematical discussion of these concepts, along with some speculative ideas, is given in "The Cosmic Origins of Time's Arrow," Sean M. Carroll, Scientific American, June 2008, p. 48.

5.4.6 Quantum mechanics and zero entropy

The previous discussion would seem to imply that absolute entropies are never well defined, since any calculation of entropy will always end up having terms that depend on Δp and Δx . For instance, we might think that cooling an ideal gas to absolute zero would give zero entropy, since there is then only one available momentum state, but there would still be many possible position states. We'll see later in this book, however, that the quantum mechanical uncertainty principle makes it impossible to know the location and position of a particle simultaneously with perfect accuracy. The best we can do is to determine them with an accuracy such that the product $\Delta p\Delta x$ is equal to a constant called Planck's constant. According to quantum physics, then, there is a natural minimum size for rectangles in phase space, and entropy can be defined in absolute terms. Another way of looking at it is that according to quantum physics, the gas as a whole has some well-defined ground state, which is its state of minimum energy. When the gas is cooled to absolute zero, the scene is not at all like what we would picture in classical physics, with a lot of atoms lying around motionless. It might, for instance, be a strange quantum-mechanical state called the Bose-Einstein condensate, which was achieved for the first time recently with macroscopic amounts of atoms. Classically, the gas has many possible states available to it at zero temperature, since the positions of the atoms can be chosen in a variety of ways. The classical picture is a bad approximation under these circumstances, however. Quantum mechanically there is only one ground state, in which each atom is spread out over the available volume in a cloud of probability. The entropy is therefore zero at zero temperature. This fact, which cannot be understood in terms of classical physics, is known as the third law of thermodynamics.

5.4.7 Summary of the laws of thermodynamics

Here is a summary of the laws of thermodynamics:

The zeroth law of thermodynamics (page 313) If object A is at the same temperature as object B, and B is at the same temperature as C, then A is at the same temperature as C.

The first law of thermodynamics (page 308) Energy is conserved.

The second law of thermodynamics (page 324) The entropy of a closed system always increases, or at best stays the same: $\Delta S \geq 0$.

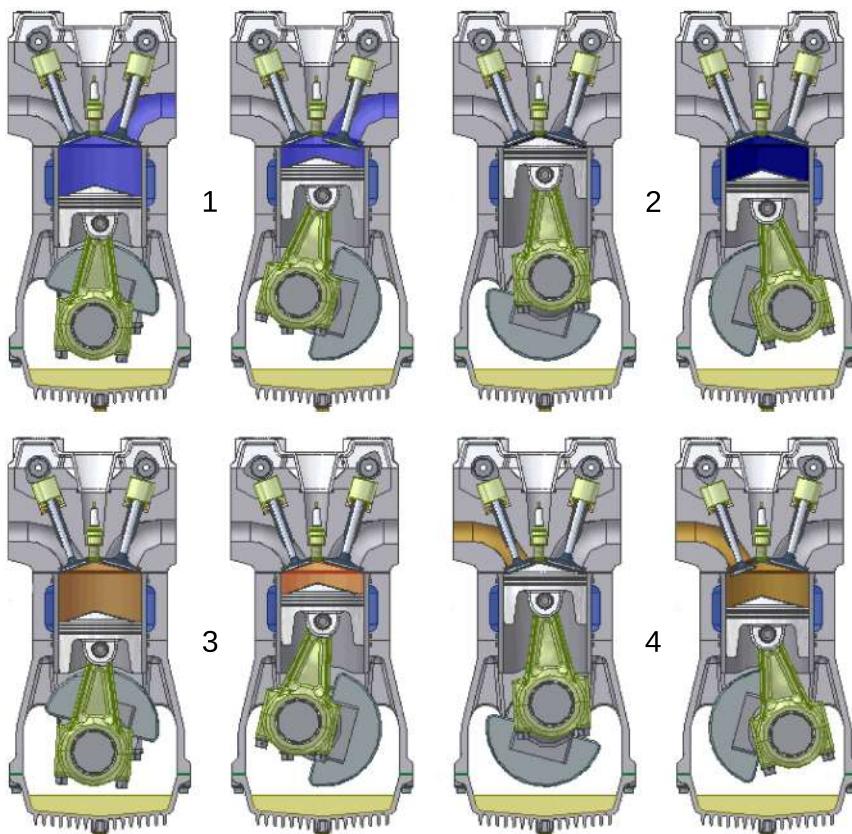
The third law of thermodynamics (page 339) The entropy of a system approaches zero as its temperature approaches absolute zero.

From a modern point of view, only the first law deserves to be called a fundamental law of physics. Once Boltzmann discovered the microscopic nature of entropy, the zeroth and second laws could be understood as statements about probability: a system containing a large number of particles is overwhelmingly likely to do a certain thing, simply because the number of possible ways to do it is extremely large compared to the other possibilities. The third law is also now understood to be a consequence of more basic physical principles, but to explain the third law, it's not sufficient simply to know that matter is made of atoms: we also need to understand the quantum-mechanical nature of those atoms, discussed in chapter 13. Historically, however, the laws of thermodynamics were discovered in the eighteenth century, when the atomic theory of matter was generally considered to be a hypothesis that couldn't be tested experimentally. Ideally, with the publication of Boltzmann's work on entropy in 1877, the zeroth and second laws would have been immediately demoted from the status of physical laws, and likewise the development of quantum mechanics in the 1920's would have done the same for the third law.

5.5 More about heat engines

So far, the only heat engine we've discussed in any detail has been a fictitious Carnot engine, with a monoatomic ideal gas as its working gas. As a more realistic example, figure m shows one full cycle of a cylinder in a standard gas-burning automobile engine. This four-stroke cycle is called the Otto cycle, after its inventor, German engineer Nikolaus Otto. The Otto cycle is more complicated than a Carnot cycle, in a number of ways:

- The working gas is physically pumped in and out of the cylinder through valves, rather than being sealed and reused indefinitely as in the Carnot engine.
- The cylinders are not perfectly insulated from the engine block, so heat energy is lost from each cylinder by conduction. This makes the engine less efficient than a Carnot engine, because heat is being discharged at a temperature that is not as cool as the environment.



m / The Otto cycle. 1. In the exhaust stroke, the piston expels the burned air-gas mixture left over from the preceding cycle. 2. In the intake stroke, the piston sucks in fresh air-gas mixture. 3. In the compression stroke, the piston compresses the mixture, and heats it. 4. At the beginning of the power stroke, the spark plug fires, causing the air-gas mixture to burn explosively and heat up much more. The heated mixture expands, and does a large amount of positive mechanical work on the piston. An animated version can be viewed in the Wikipedia article “Four-stroke engine.”

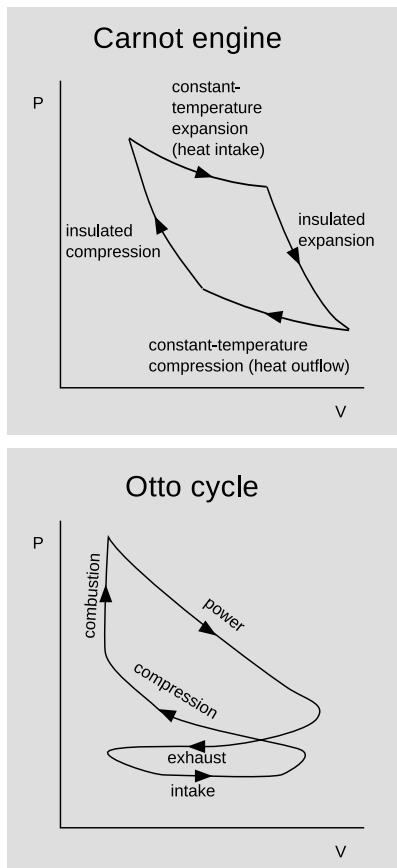
- Rather than being heated by contact with an external heat reservoir, the air-gas mixture inside each cylinder is heated by internal combustion: a spark from a spark plug burns the gasoline, releasing heat.
- The working gas is not monoatomic. Air consists of diatomic molecules (N_2 and O_2), and gasoline of polyatomic molecules such as octane (C_8H_{18}).
- The working gas is not ideal. An ideal gas is one in which the molecules never interact with one another, but only with the walls of the vessel, when they collide with it. In a car engine, the molecules are interacting very dramatically with one another when the air-gas mixture explodes (and less dramatically at other times as well, since, for example, the gasoline may be in the form of microscopic droplets rather than individual molecules).

This is all extremely complicated, and it would be nice to have some way of understanding and visualizing the important properties of such a heat engine without trying to handle every detail at once. A good method of doing this is a type of graph known as a P-V diagram. As proved in homework problem 2, the equation $dW = F dx$ for mechanical work can be rewritten as $dW = P dV$ in

the case of work done by a piston. Here P represents the pressure of the working gas, and V its volume. Thus, on a graph of P versus V , the area under the curve represents the work done. When the gas expands, dx is positive, and the gas does positive work. When the gas is being compressed, dx is negative, and the gas does negative work, i.e., it absorbs energy. Notice how, in the diagram of the Carnot engine in the top panel of figure a, the cycle goes clockwise around the curve, and therefore the part of the curve in which negative work is being done (arrowheads pointing to the left) are below the ones in which positive work is being done. This means that over all, the engine does a positive amount of work. This net work equals the area under the top part of the curve, minus the area under the bottom part of the curve, which is simply the area enclosed by the curve. Although the diagram for the Otto engine is more complicated, we can at least compare it on the same footing with the Carnot engine. The curve forms a figure-eight, because it cuts across itself. The top loop goes clockwise, so as in the case of the Carnot engine, it represents positive work. The bottom loop goes counterclockwise, so it represents a net negative contribution to the work. This is because more work is expended in forcing out the exhaust than is generated in the intake stroke.

To make an engine as efficient as possible, we would like to make the loop have as much area as possible. What is it that determines the actual shape of the curve? First let's consider the constant-temperature expansion stroke that forms the top of the Carnot engine's P-V plot. This is analogous to the power stroke of an Otto engine. Heat is being sucked in from the hot reservoir, and since the working gas is always in thermal equilibrium with the hot reservoir, its temperature is constant. Regardless of the type of gas, we therefore have $PV = nkT$ with T held constant, and thus $P \propto V^{-1}$ is the mathematical shape of this curve — a $y = 1/x$ graph, which is a hyperbola. This is all true regardless of whether the working gas is monoatomic, diatomic, or polyatomic. (The bottom of the loop is likewise of the form $P \propto V^{-1}$, but with a smaller constant of proportionality due to the lower temperature.)

Now consider the insulated expansion stroke that forms the right side of the curve for the Carnot engine. As shown on page 336, the relationship between pressure and temperature in an insulated compression or expansion is $T \propto P^b$, with $b = 2/5, 2/7$, or $1/4$, respectively, for a monoatomic, diatomic, or polyatomic gas. For P as a function of V at constant T , the ideal gas law gives $P \propto T/V$, so $P \propto V^{-\gamma}$, where $\gamma = 1/(1 - b)$ takes on the values $5/3, 7/5$, and $4/3$. The number γ can be interpreted as the ratio C_P/C_V , where C_P , the heat capacity at constant pressure, is the amount of heat required to raise the temperature of the gas by one degree while keeping its pressure constant, and C_V is the corresponding quantity under conditions of constant volume.



a / P-V diagrams for a Carnot engine and an Otto engine.

The compression ratio

example 22

Operating along a constant-temperature stroke, the amount of mechanical work done by a heat engine can be calculated as follows:

$$PV = nkT$$

Setting $c = nkT$ to simplify the writing,

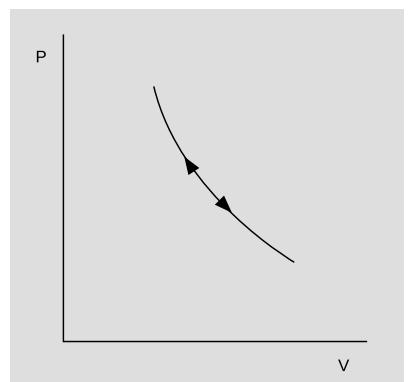
$$\begin{aligned} P &= cV^{-1} \\ W &= \int_{V_i}^{V_f} P dV \\ &= c \int_{V_i}^{V_f} V^{-1} dV \\ &= c \ln V_f - c \ln V_i \\ &= c \ln(V_f/V_i) \end{aligned}$$

The ratio V_f/V_i is called the *compression ratio* of the engine, and higher values result in more power along this stroke. Along an insulated stroke, we have $P \propto V^{-\gamma}$, with $\gamma \neq 1$, so the result for the work no longer has this perfect mathematical property of depending only on the ratio V_f/V_i . Nevertheless, the compression ratio is still a good figure of merit for predicting the performance of any heat engine, including an internal combustion engine. High compression ratios tend to make the working gas of an internal combustion engine heat up so much that it spontaneously explodes. When this happens in an Otto-cycle engine, it can cause ignition before the sparkplug fires, an undesirable effect known as pinging. For this reason, the compression ratio of an Otto-cycle automobile engine cannot normally exceed about 10. In a diesel engine, however, this effect is used intentionally, as an alternative to sparkplugs, and compression ratios can be 20 or more.

Sound

example 23

Figure b shows a P-V plot for a sound wave. As the pressure oscillates up and down, the air is heated and cooled by its compression and expansion. Heat conduction is a relatively slow process, so typically there is not enough time over each cycle for any significant amount of heat to flow from the hot areas to the cold areas. (This is analogous to insulated compression or expansion of a heat engine; in general, a compression or expansion of this type, with no transfer of heat, is called *adiabatic*.) The pressure and volume of a particular little piece of the air are therefore related according to $P \propto V^{-\gamma}$. The cycle of oscillation consists of motion back and forth along a single curve in the P-V plane, and since this curve encloses zero volume, no mechanical work



b / Example 23,

is being done: the wave (under the assumed ideal conditions) propagates without any loss of energy due to friction.

The speed of sound is also related to γ . See example 13 on p. 389.

Measuring γ using the “spring of air”

example 24

Figure c shows an experiment that can be used to measure the γ of a gas. When the mass m is inserted into bottle's neck, which has cross-sectional area A , the mass drops until it compresses the air enough so that the pressure is enough to support its weight. The observed frequency ω of oscillations about this equilibrium position y_0 can be used to extract the γ of the gas.

$$\begin{aligned}\omega^2 &= \frac{k}{m} \\ &= -\frac{1}{m} \left. \frac{dF}{dy} \right|_{y_0} \\ &= -\left. \frac{A}{m} \frac{dP}{dy} \right|_{y_0} \\ &= -\left. \frac{A^2}{m} \frac{dP}{dV} \right|_{V_0}\end{aligned}$$

c / Example 24.

We make the bottle big enough so that its large surface-to-volume ratio prevents the conduction of any significant amount of heat through its walls during one cycle, so $P \propto V^{-\gamma}$, and $dP/dV = -\gamma P/V$. Thus,

$$\omega^2 = \gamma \frac{A^2}{m} \frac{P_0}{V_0}$$

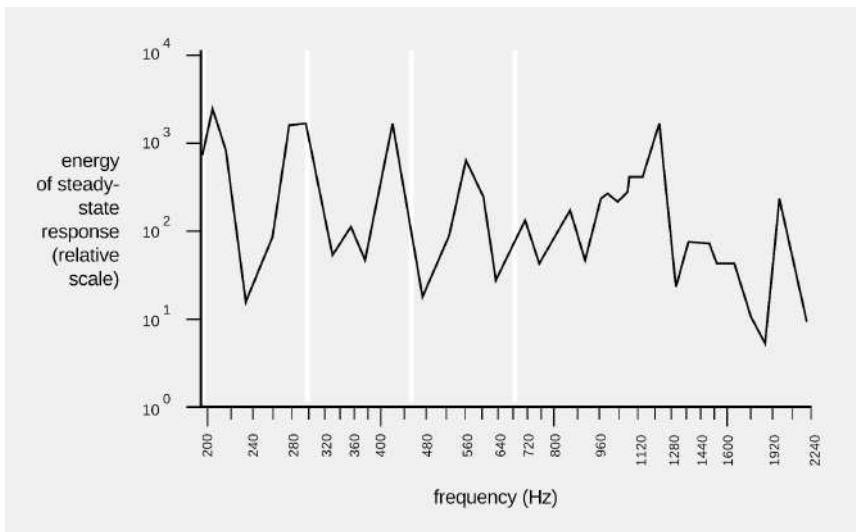
The Helmholtz resonator

example 25

When you blow over the top of a beer bottle, you produce a pure tone. As you drink more of the beer, the pitch goes down. This is similar to example 24, except that instead of a solid mass m sitting inside the neck of the bottle, the moving mass is the air itself. As air rushes in and out of the bottle, its velocity is highest at the bottleneck, and since kinetic energy is proportional to the square of the velocity, essentially all of the kinetic energy is that of the air that's in the neck. In other words, we can replace m with $AL\rho$, where L is the length of the neck, and ρ is the density of the air. Substituting into the earlier result, we find that the resonant frequency is

$$\omega^2 = \gamma \frac{P_0}{\rho} \frac{A}{LV_0}.$$

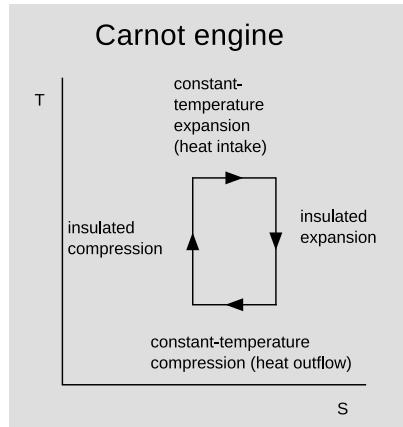
This is known as a Helmholtz resonator. As shown in figure d, a violin or an acoustic guitar has a Helmholtz resonance, since air



d / The resonance curve of a 1713 Stradivarius violin, measured by Carleen Hutchins. There are a number of different resonance peaks, some strong and some weak; the ones near 200 and 400 Hz are vibrations of the wood, but the one near 300 Hz is a resonance of the air moving in and out through those holes shaped like the letter F. The white lines show the frequencies of the four strings.

can move in and out through the f-holes. Problem 10 is a more quantitative exploration of this.

We have already seen, based on the microscopic nature of entropy, that any Carnot engine has the same efficiency, and the argument only employed the assumption that the engine met the definition of a Carnot cycle: two insulated strokes, and two constant-temperature strokes. Since we didn't have to make any assumptions about the nature of the working gas being used, the result is evidently true for diatomic or polyatomic molecules, or for a gas that is not ideal. This result is surprisingly simple and general, and a little mysterious — it even applies to possibilities that we have not even considered, such as a Carnot engine designed so that the working “gas” actually consists of a mixture of liquid droplets and vapor, as in a steam engine. How can it always turn out so simple, given the kind of mathematical complications that were swept under the rug in example 22? A better way to understand this result is by switching from P-V diagrams to a diagram of temperature versus entropy, as shown in figure e. An infinitesimal transfer of heat dQ gives rise to a change in entropy $dS = dQ/T$, so the area under the curve on a T-S plot gives the amount of heat transferred. The area under the top edge of the box in figure e, extending all the way down to the axis, represents the amount of heat absorbed from the hot reservoir, while the smaller area under the bottom edge represents the heat wasted into the cold reservoir. By conservation of energy, the area enclosed by the box therefore represents the amount of mechanical work being done, as for a P-V diagram. We can now see why the efficiency of a Carnot engine is independent of any of the physical details: the definition of a Carnot engine guarantees that the T-S diagram will be a rectangular box, and the efficiency depends only on the relative heights of the top and bottom of the box.



e / A T-S diagram for a Carnot engine.

This chapter is summarized on page 1081. Notation and terminology are tabulated on pages 1070-1071.

Problems

The symbols \checkmark , \blacksquare , etc. are explained on page 352.

1 (a) Show that under conditions of standard pressure and temperature, the volume of a sample of an ideal gas depends only on the number of molecules in it.

(b) One mole is defined as 6.0×10^{23} atoms. Find the volume of one mole of an ideal gas, in units of liters, at standard temperature and pressure (0°C and 101 kPa). \checkmark \blacksquare

2 A gas in a cylinder expands its volume by an amount dV , pushing out a piston. Show that the work done by the gas on the piston is given by $dW = P dV$. \blacksquare

3 (a) A helium atom contains 2 protons, 2 electrons, and 2 neutrons. Find the mass of a helium atom. \checkmark

(b) Find the number of atoms in 1.0 kg of helium. \checkmark

(c) Helium gas is monoatomic. Find the amount of heat needed to raise the temperature of 1.0 kg of helium by 1.0 degree C. (This is known as helium's heat capacity at constant volume.) \checkmark \blacksquare

4 A sample of gas is enclosed in a sealed chamber. The gas consists of molecules, which are then split in half through some process such as exposure to ultraviolet light, or passing an electric spark through the gas. The gas returns to the same temperature as the surrounding room, but the molecules remain split apart, at least for some amount of time. (To achieve these conditions, we would need an extremely dilute gas. Otherwise the recombination of the molecules would be faster than the cooling down to the same temperature as the room.) How does its pressure now compare with its pressure before the molecules were split? \blacksquare

5 Most of the atoms in the universe are in the form of gas that is not part of any star or galaxy: the intergalactic medium (IGM). The IGM consists of about 10^{-5} atoms per cubic centimeter, with a typical temperature of about 10^3 K. These are, in some sense, the density and temperature of the universe (not counting light, or the exotic particles known as "dark matter"). Calculate the pressure of the universe (or, speaking more carefully, the typical pressure due to the IGM). \checkmark \blacksquare

6 Estimate the pressure at the center of the Earth, assuming it is of constant density throughout. Note that g is not constant with respect to depth — as shown in example 19 on page 105, g equals Gmr/b^3 for r , the distance from the center, less than b , the earth's radius.

(a) State your result in terms of G , m , and b . \checkmark

(b) Show that your answer from part a has the right units for pressure.

(c) Evaluate the result numerically. \checkmark

(d) Given that the earth's atmosphere is on the order of one thou-

sandth the earth's radius, and that the density of the earth is several thousand times greater than the density of the lower atmosphere, check that your result is of a reasonable order of magnitude. ■

7 (a) Determine the ratio between the escape velocities from the surfaces of the earth and the moon. ✓

(b) The temperature during the lunar daytime gets up to about 130°C. In the extremely thin (almost nonexistent) lunar atmosphere, estimate how the typical velocity of a molecule would compare with that of the same type of molecule in the earth's atmosphere. Assume that the earth's atmosphere has a temperature of 0°C. ✓

(c) Suppose you were to go to the moon and release some fluorocarbon gas, with molecular formula C_nF_{2n+2} . Estimate what is the smallest fluorocarbon molecule (lowest n) whose typical velocity would be lower than that of an N_2 molecule on earth in proportion to the moon's lower escape velocity. The moon would be able to retain an atmosphere made of these molecules. ✓ ■

8 Refrigerators, air conditioners, and heat pumps are heat engines that work in reverse. You put in mechanical work, and the effect is to take heat out of a cooler reservoir and deposit heat in a warmer one: $Q_L + W = Q_H$. As with the heat engines discussed previously, the efficiency is defined as the energy transfer you want (Q_L for a refrigerator or air conditioner, Q_H for a heat pump) divided by the energy transfer you pay for (W).

Efficiencies are supposed to be unitless, but the efficiency of an air conditioner is normally given in terms of an EER rating (or a more complex version called an SEER). The EER is defined as Q_L/W , but expressed in the barbaric units of of Btu/watt-hour. A typical EER rating for a residential air conditioner is about 10 Btu/watt-hour, corresponding to an efficiency of about 3. The standard temperatures used for testing an air conditioner's efficiency are 80°F (27°C) inside and 95°F (35°C) outside.

(a) What would be the EER rating of a reversed Carnot engine used as an air conditioner? ✓

(b) If you ran a 3-kW residential air conditioner, with an efficiency of 3, for one hour, what would be the effect on the total entropy of the universe? Is your answer consistent with the second law of thermodynamics? ✓ ■

9 Even when resting, the human body needs to do a certain amount of mechanical work to keep the heart beating. This quantity is difficult to define and measure with high precision, and also depends on the individual and her level of activity, but it's estimated to be about 1 to 5 watts. Suppose we consider the human body as nothing more than a pump. A person who is just lying in bed all day needs about 1000 kcal/day worth of food to stay alive. (a) Estimate the person's thermodynamic efficiency as a pump, and (b) compare with the maximum possible efficiency imposed by the laws

of thermodynamics for a heat engine operating across the difference between a body temperature of 37°C and an ambient temperature of 22°C . (c) Interpret your answer. ▷ Answer, p. 1068 ■

10 Example 25 on page 344 suggests analyzing the resonance of a violin at 300 Hz as a Helmholtz resonance. However, we might expect the equation for the frequency of a Helmholtz resonator to be a rather crude approximation here, since the f-holes are not long tubes, but slits cut through the face of the instrument, which is only about 2.5 mm thick. (a) Estimate the frequency that way anyway, for a violin with a volume of about 1.6 liters, and f-holes with a total area of 10 cm^2 . (b) A common rule of thumb is that at an open end of an air column, such as the neck of a real Helmholtz resonator, some air beyond the mouth also vibrates as if it was inside the tube, and that this effect can be taken into account by adding 0.4 times the diameter of the tube for each open end (i.e., 0.8 times the diameter when both ends are open). Applying this to the violin's f-holes results in a huge change in L , since the $\sim 7 \text{ mm}$ width of the f-hole is considerably greater than the thickness of the wood. Try it, and see if the result is a better approximation to the observed frequency of the resonance. ▷ Answer, p. 1068 ■

11 (a) Atmospheric pressure at sea level is 101 kPa. The deepest spot in the world's oceans is a valley called the Challenger Deep, in the Marianas Trench, with a depth of 11.0 km. Find the pressure at this depth, in units of atmospheres. Although water under this amount of pressure does compress by a few percent, assume for the purposes of this problem that it is incompressible.

(b) Suppose that an air bubble is formed at this depth and then rises to the surface. Estimate the change in its volume and radius.

▷ Solution, p. 1044 ■

12 Our sun is powered by nuclear fusion reactions, and as a first step in these reactions, one proton must approach another proton to within a short enough range r . This is difficult to achieve, because the protons have electric charge $+e$ and therefore repel one another electrically. (It's a good thing that it's so difficult, because otherwise the sun would use up all of its fuel very rapidly and explode.) To make fusion possible, the protons must be moving fast enough to come within the required range. Even at the high temperatures present in the core of our sun, almost none of the protons are moving fast enough.

(a) For comparison, the early universe, soon after the Big Bang, had extremely high temperatures. Estimate the temperature T that would have been required so that protons with average energies could fuse. State your result in terms of r , the mass m of the proton, and universal constants.

(b) Show that the units of your answer to part a make sense.

(c) Evaluate your result from part a numerically, using $r = 10^{-15} \text{ m}$ and $m = 1.7 \times 10^{-27} \text{ kg}$. As a check, you should find that this is

much hotter than the sun's core temperature of $\sim 10^7$ K.

▷ Solution, p. 1045 ■

13 Object A is a brick. Object B is half of a similar brick. If A is heated, we have $\Delta S = Q/T$. Show that if this equation is valid for A, then it is also valid for B. \triangleright Solution, p. 1045 ■

14 Typically the atmosphere gets colder with increasing altitude. However, sometimes there is an *inversion layer*, in which this trend is reversed, e.g., because a less dense mass of warm air moves into a certain area, and rises above the denser colder air that was already present. Suppose that this causes the pressure P as a function of height y to be given by a function of the form $P = P_o e^{-ky}(1 + by)$, where constant temperature would give $b = 0$ and an inversion layer would give $b > 0$. (a) Infer the units of the constants P_o , k , and b . (b) Find the density of the air as a function of y , of the constants, and of the acceleration of gravity g . (c) Check that the units of your answer to part b make sense. \triangleright Solution, p. 1045 ■

15 (a) Consider a *one-dimensional* ideal gas consisting of n material particles, at temperature T . Trace back through the logic of the equipartition theorem on p. 333 to determine the total energy. (b) Explain why it should matter how many dimensions there are. (c) Gases that we encounter in everyday life are made of atoms, but there are gases made out of other things. For example, soon after the big bang, there was a period when the universe was very hot and dominated by light rather than matter. A particle of light is called a photon, so the early universe was a “photon gas.” For simplicity, consider a photon gas in one dimension. Photons are massless, and we will see in ch. 7 on relativity that for a massless particle, the energy is related to the momentum by $E = pc$, where c is the speed of light. (Note that $p = mv$ does *not* hold for a photon.) Again, trace back through the logic of equipartition on p. 333. Does the photon gas have the same heat capacity as the one you found in part a? ■

16 You use a spoon at room temperature, 22°C , to mix your coffee, which is at 80°C . During this brief period of thermal contact, 1.3 J of heat is transferred from the coffee to the spoon. Find the total change in the entropy of the universe. \checkmark ■

17 The sun is mainly a mixture of hydrogen and helium, some of which is ionized. As a simplified model, let's pretend that it's made purely out of neutral, monatomic hydrogen, and that the whole mass of the sun is in thermal equilibrium. Given its mass, it would then contain 1.2×10^{57} atoms. It generates energy from nuclear reactions at a rate of $3.8 \times 10^{26} \text{ W}$, and it is in a state of equilibrium in which this amount of energy is radiated off into space as light. Suppose that its ability to radiate light were somehow blocked. Find the rate at which its temperature would increase. \checkmark ■

18 In metals, some electrons, called conduction electrons, are free to move around, rather than being bound to one atom. Classical physics gives an adequate description of many of their properties. Consider a metal at temperature T , and let m be the mass of the electron. Find expressions for (a) the average kinetic energy of a conduction electron, and (b) the average square of its velocity, v^2 . (It would not be of much interest to find \bar{v} , which is just zero.) Numerically, $\sqrt{\bar{v}^2}$, called the root-mean-square velocity, comes out to be surprisingly large — about two orders of magnitude greater than the normal thermal velocities we find for atoms in a gas. Why?

Remark: From this analysis, one would think that the conduction electrons would contribute greatly to the heat capacities of metals. In fact they do not contribute very much in most cases; if they did, Dulong and Petit's observations would not have come out as described in the text. The resolution of this contradiction was only eventually worked out by Sommerfeld in 1933, and involves the fact that electrons obey the Pauli exclusion principle. \checkmark ■

Key to symbols:

■ easy ■ typical ■ challenging ■ difficult ■ very difficult
 \checkmark An answer check is available at www.lightandmatter.com.



The vibrations of this electric bass string are converted to electrical vibrations, then to sound vibrations, and finally to vibrations of our eardrums.

Chapter 6

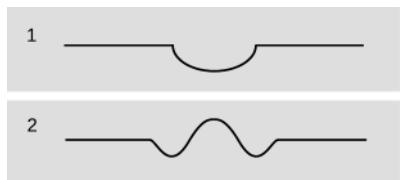
Waves

Dandelion. Cello. Read those two words, and your brain instantly conjures a stream of associations, the most prominent of which have to do with vibrations. Our mental category of “dandelion-ness” is strongly linked to the color of light waves that vibrate about half a million billion times a second: yellow. The velvety throb of a cello has as its most obvious characteristic a relatively low musical pitch — the note you’re spontaneously imagining right now might be one whose sound vibrations repeat at a rate of a hundred times a second.

Evolution seems to have designed our two most important senses around the assumption that our environment is made of waves, whereas up until now, we’ve mostly taken the view that Nature can be understood by breaking her down into smaller and smaller parts, ending up with particles as her most fundamental building blocks. Does that work for light and sound? Sound waves are disturbances in air, which is made of atoms, but light, on the other hand, isn’t a vibration of atoms. Light, unlike sound, can travel through a vacuum: if you’re reading this by sunlight, you’re taking advantage of light that had to make it through millions of miles of vacuum to get to you. Waves, then, are not just a trick that vibrating atoms can do. Waves are one of the basic phenomena of the universe. At the

end of this book, we'll even see that the things we've been calling particles, such as electrons, are really waves!¹

6.1 Free waves



a / Your finger makes a depression in the surface of the water, 1. The wave pattern starts evolving, 2, after you remove your finger.

6.1.1 Wave motion

Let's start with an intuition-building exercise that deals with waves in matter, since they're easier than light waves to get your hands on. Put your fingertip in the middle of a cup of water and then remove it suddenly. You'll have noticed two results that are surprising to most people. First, the flat surface of the water does not simply sink uniformly to fill in the volume vacated by your finger. Instead, ripples spread out, and the process of flattening out occurs over a long period of time, during which the water at the center vibrates above and below the normal water level. This type of wave motion is the topic of the present section. Second, you've found that the ripples bounce off of the walls of the cup, in much the same way that a ball would bounce off of a wall. In the next section we discuss what happens to waves that have a boundary around them. Until then, we confine ourselves to wave phenomena that can be analyzed as if the medium (e.g., the water) was infinite and the same everywhere.

It isn't hard to understand why removing your fingertip creates ripples rather than simply allowing the water to sink back down uniformly. The initial crater, a/1, left behind by your finger has sloping sides, and the water next to the crater flows downhill to fill in the hole. The water far away, on the other hand, initially has no way of knowing what has happened, because there is no slope for it to flow down. As the hole fills up, the rising water at the center gains upward momentum, and overshoots, creating a little hill where there had been a hole originally. The area just outside of this region has been robbed of some of its water in order to build the hill, so a depressed "moat" is formed, a/2. This effect cascades outward, producing ripples.

There are three main ways in which wave motion differs from the motion of objects made of matter.

1. Superposition

If you watched the water in the cup carefully, you noticed the ghostlike behavior of the reflected ripples coming back toward the center of the cup and the outgoing ripples that hadn't yet been reflected: they passed right through each other. This is the first, and the most profound, difference between wave motion and the mo-

¹Speaking more carefully, I should say that every basic building block of light and matter has both wave and particle properties.

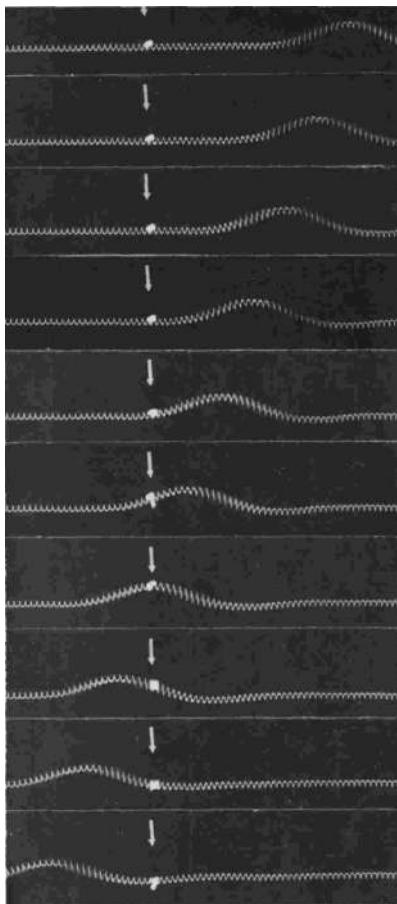


b / The two circular patterns of ripples pass through each other. Unlike material objects, wave patterns can overlap in space, and when this happens they combine by addition.

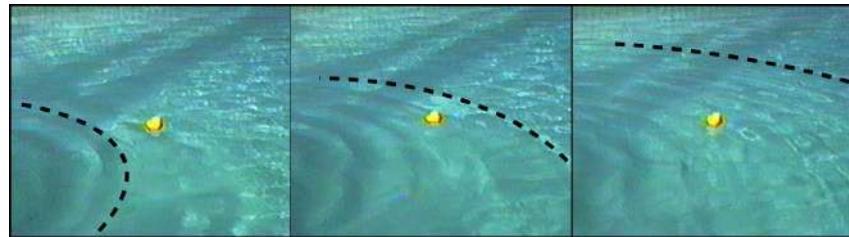
tion of objects: waves do not display any repulsion of each other analogous to the normal forces between objects that come in contact. Two wave patterns can therefore overlap in the same region of space, as shown in figure b. Where the two waves coincide, they add together. For instance, suppose that at a certain location in at a certain moment in time, each wave would have had a crest 3 cm above the normal water level. The waves combine at this point to make a 6-cm crest. We use negative numbers to represent depressions in the water. If both waves would have had a troughs measuring -3 cm, then they combine to make an extra-deep -6 cm trough. A $+3$ cm crest and a -3 cm trough result in a height of zero, i.e., the waves momentarily cancel each other out at that point. This additive rule is referred to as the principle of superposition, “superposition” being merely a fancy word for “adding.”

Superposition can occur not just with sinusoidal waves like the ones in the figure above but with waves of any shape. The figures on the following page show superposition of wave pulses. A pulse is simply a wave of very short duration. These pulses consist only of a single hump or trough. If you hit a clothesline sharply, you will observe pulses heading off in both directions. This is analogous to the way ripples spread out in all directions when you make a disturbance at one point on water. The same occurs when the hammer on a piano comes up and hits a string.

Experiments to date have not shown any deviation from the principle of superposition in the case of light waves. For other types of waves, it is typically a very good approximation for low-energy waves.



d / As the wave pulse goes by, the ribbon tied to the spring is not carried along. The motion of the wave pattern is to the right, but the medium (spring) is moving from side to side, not to the right. (*PSSC Physics*)



c / As the wave pattern passes the rubber duck, the duck stays put. The water isn't moving with the wave.

2. The medium is not transported with the wave.

The sequence of three photos in figure c shows a series of water waves before it has reached a rubber duck (left), having just passed the duck (middle) and having progressed about a meter beyond the duck (right). The duck bobs around its initial position, but is not carried along with the wave. This shows that the water itself does not flow outward with the wave. If it did, we could empty one end of a swimming pool simply by kicking up waves! We must distinguish between the motion of the medium (water in this case) and the motion of the wave pattern through the medium. The medium vibrates; the wave progresses through space.

self-check A

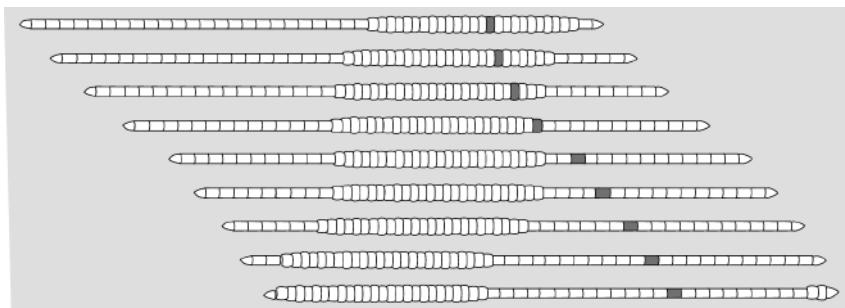
In figure d, you can detect the side-to-side motion of the spring because the spring appears blurry. At a certain instant, represented by a single photo, how would you describe the motion of the different parts of the spring? Other than the flat parts, do any parts of the spring have zero velocity?

▷ Answer, p. 1061

A worm

example 1

The worm in the figure is moving to the right. The wave pattern, a pulse consisting of a compressed area of its body, moves to the left. In other words, the motion of the wave pattern is in the opposite direction compared to the motion of the medium.



Surfing

example 2

The incorrect belief that the medium moves with the wave is often reinforced by garbled secondhand knowledge of surfing. Anyone who has actually surfed knows that the front of the board pushes the water to the sides, creating a wake — the surfer can even drag his hand through the water, as in figure e. If the water was moving along with the wave and the surfer, this wouldn't happen. The surfer is carried forward because forward is downhill, not because of any forward flow of the water. If the water was flowing forward, then a person floating in the water up to her neck would be carried along just as quickly as someone on a surfboard. In fact, it is even possible to surf down the back side of a wave, although the ride wouldn't last very long because the surfer and the wave would quickly part company.

3. A wave's velocity depends on the medium.

A material object can move with any velocity, and can be sped up or slowed down by a force that increases or decreases its kinetic energy. Not so with waves. The speed of a wave, depends on the properties of the medium (and perhaps also on the shape of the wave, for certain types of waves). Sound waves travel at about 340 m/s in air, 1000 m/s in helium. If you kick up water waves in a pool, you will find that kicking harder makes waves that are taller (and therefore carry more energy), not faster. The sound waves from an exploding stick of dynamite carry a lot of energy, but are no faster than any other waves. In the following section we will give an example of the physical relationship between the wave speed and the properties of the medium.

Breaking waves

example 3

The velocity of water waves increases with depth. The crest of a wave travels faster than the trough, and this can cause the wave to break.

Once a wave is created, the only reason its speed will change is if it enters a different medium or if the properties of the medium change. It is not so surprising that a change in medium can slow down a wave, but the reverse can also happen. A sound wave traveling through a helium balloon will slow down when it emerges into the air, but if it enters another balloon it will speed back up again! Similarly, water waves travel more quickly over deeper water, so a wave will slow down as it passes over an underwater ridge, but speed up again as it emerges into deeper water.

Hull speed

example 4

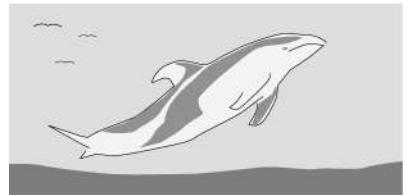
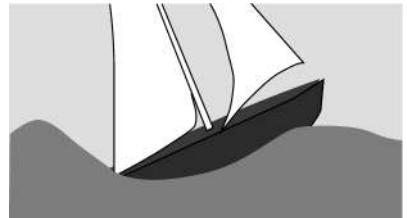
The speeds of most boats, and of some surface-swimming animals, are limited by the fact that they make a wave due to their motion through the water. The boat in figure g is going at the same speed as its own waves, and can't go any faster. No matter how hard the boat pushes against the water, it can't make



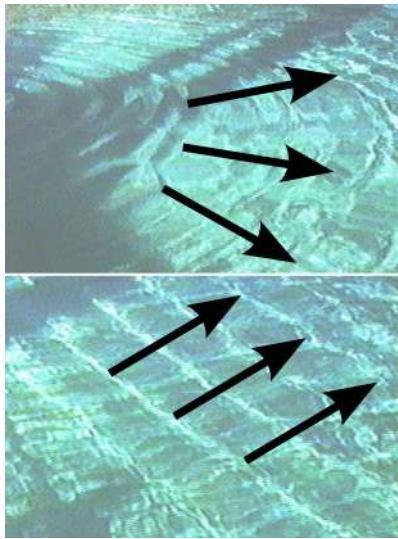
e / Example 2. The surfer is dragging his hand in the water.



f / Example 3: a breaking wave.



g / Example 4. The boat has run up against a limit on its speed because it can't climb over its own wave. Dolphins get around the problem by leaping out of the water.



h / Circular and linear wave patterns.

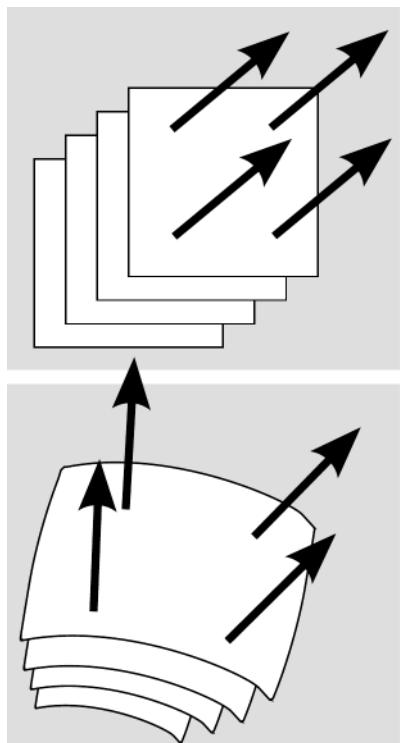
the wave move ahead faster and get out of the way. The wave's speed depends only on the medium. Adding energy to the wave doesn't speed it up, it just increases its amplitude.

A water wave, unlike many other types of wave, has a speed that depends on its shape: a broader wave moves faster. The shape of the wave made by a boat tends to mold itself to the shape of the boat's hull, so a boat with a longer hull makes a broader wave that moves faster. The maximum speed of a boat whose speed is limited by this effect is therefore closely related to the length of its hull, and the maximum speed is called the hull speed. Sailboats designed for racing are not just long and skinny to make them more streamlined — they are also long so that their hull speeds will be high.

Wave patterns

If the magnitude of a wave's velocity vector is preordained, what about its direction? Waves spread out in all directions from every point on the disturbance that created them. If the disturbance is small, we may consider it as a single point, and in the case of water waves the resulting wave pattern is the familiar circular ripple, *h/1*. If, on the other hand, we lay a pole on the surface of the water and wiggle it up and down, we create a linear wave pattern, *h/2*. For a three-dimensional wave such as a sound wave, the analogous patterns would be spherical waves and plane waves, *i*.

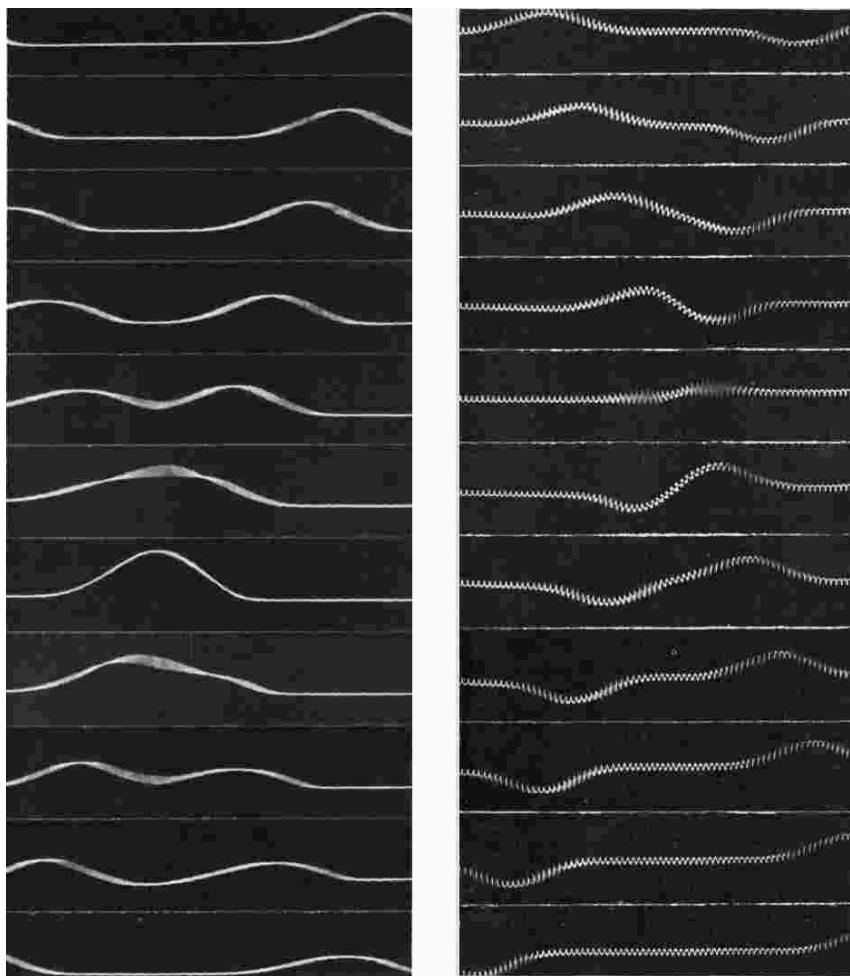
Infinitely many patterns are possible, but linear or plane waves are often the simplest to analyze, because the velocity vector is in the same direction no matter what part of the wave we look at. Since all the velocity vectors are parallel to one another, the problem is effectively one-dimensional. Throughout this chapter and the next, we will restrict ourselves mainly to wave motion in one dimension, while not hesitating to broaden our horizons when it can be done without too much complication.



i / Plane and spherical wave patterns.

Discussion Questions

A The left panel of the figure shows a sequence of snapshots of two positive pulses on a coil spring as they move through each other. In the right panel, which shows a positive pulse and a negative one, the fifth frame has the spring just about perfectly flat. If the two pulses have essentially canceled each other out perfectly, then why does the motion pick up again? Why doesn't the spring just stay flat?



j / Discussion question A.

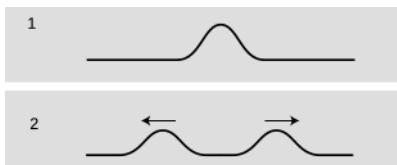
B Sketch two positive wave pulses on a string that are overlapping but not right on top of each other, and draw their superposition. Do the same for a positive pulse running into a negative pulse.

C A traveling wave pulse is moving to the right on a string. Sketch the velocity vectors of the various parts of the string. Now do the same for a pulse moving to the left.

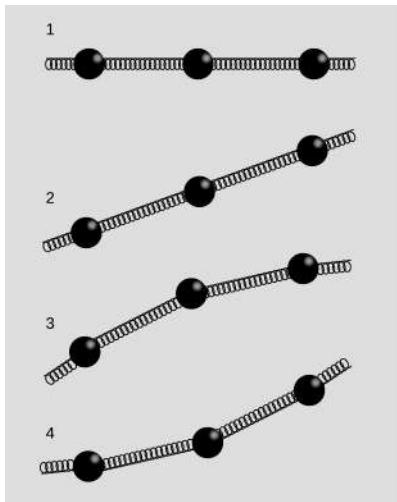
D In a spherical sound wave spreading out from a point, how would the energy of the wave fall off with distance?



k / Hitting a key on a piano causes a hammer to come up from underneath and hit a string (actually a set of three). The result is a pair of pulses moving away from the point of impact.



l / A pulse on a string splits in two and heads off in both directions.



m / Modeling a string as a series of masses connected by springs.

6.1.2 Waves on a string

So far you've learned some counterintuitive things about the behavior of waves, but intuition can be trained. The first half of this subsection aims to build your intuition by investigating a simple, one-dimensional type of wave: a wave on a string. If you have ever stretched a string between the bottoms of two open-mouthed cans to talk to a friend, you were putting this type of wave to work. Stringed instruments are another good example. Although we usually think of a piano wire simply as vibrating, the hammer actually strikes it quickly and makes a dent in it, which then ripples out in both directions. Since this chapter is about free waves, not bounded ones, we pretend that our string is infinitely long.

After the qualitative discussion, we will use simple approximations to investigate the speed of a wave pulse on a string. This quick and dirty treatment is then followed by a rigorous attack using the methods of calculus, which turns out to be both simpler and more general.

Intuitive ideas

Consider a string that has been struck, l/1, resulting in the creation of two wave pulses, l/2, one traveling to the left and one to the right. This is analogous to the way ripples spread out in all directions from a splash in water, but on a one-dimensional string, "all directions" becomes "both directions."

We can gain insight by modeling the string as a series of masses connected by springs, m. (In the actual string the mass and the springiness are both contributed by the molecules themselves.) If we look at various microscopic portions of the string, there will be some areas that are flat, 1, some that are sloping but not curved, 2, and some that are curved, 3 and 4. In example 1 it is clear that both the forces on the central mass cancel out, so it will not accelerate. The same is true of 2, however. Only in curved regions such as 3 and 4 is an acceleration produced. In these examples, the vector sum of the two forces acting on the central mass is not zero. The important concept is that curvature makes force: the curved areas of a wave tend to experience forces resulting in an acceleration toward the mouth of the curve. Note, however, that an uncurved portion of the string need not remain motionless. It may move at constant velocity to either side.

Approximate treatment

We now carry out an approximate treatment of the speed at which two pulses will spread out from an initial indentation on a string. For simplicity, we imagine a hammer blow that creates a triangular dent, n/1. We will estimate the amount of time, t, required until each of the pulses has traveled a distance equal to the width of the pulse itself. The velocity of the pulses is then $\pm w/t$.

As always, the velocity of a wave depends on the properties of the medium, in this case the string. The properties of the string can be summarized by two variables: the tension, T , and the mass per unit length, μ (Greek letter mu).

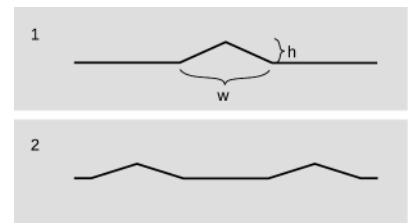
If we consider the part of the string encompassed by the initial dent as a single object, then this object has a mass of approximately μw (mass/length \times length=mass). (Here, and throughout the derivation, we assume that h is much less than w , so that we can ignore the fact that this segment of the string has a length slightly greater than w .) Although the downward acceleration of this segment of the string will be neither constant over time nor uniform across the pulse, we will pretend that it is constant for the sake of our simple estimate. Roughly speaking, the time interval between $n/1$ and $n/2$ is the amount of time required for the initial dent to accelerate from rest and reach its normal, flattened position. Of course the tip of the triangle has a longer distance to travel than the edges, but again we ignore the complications and simply assume that the segment as a whole must travel a distance h . Indeed, it might seem surprising that the triangle would so neatly spring back to a perfectly flat shape. It is an experimental fact that it does, but our analysis is too crude to address such details.

The string is kinked, i.e., tightly curved, at the edges of the triangle, so it is here that there will be large forces that do not cancel out to zero. There are two forces acting on the triangular hump, one of magnitude T acting down and to the right, and one of the same magnitude acting down and to the left. If the angle of the sloping sides is θ , then the total force on the segment equals $2T \sin \theta$. Dividing the triangle into two right triangles, we see that $\sin \theta$ equals h divided by the length of one of the sloping sides. Since h is much less than w , the length of the sloping side is essentially the same as $w/2$, so we have $\sin \theta = 2h/w$, and $F = 4Th/w$. The acceleration of the segment (actually the acceleration of its center of mass) is

$$\begin{aligned} a &= \frac{F}{m} \\ &= \frac{4Th}{\mu w^2}. \end{aligned}$$

The time required to move a distance h under constant acceleration a is found by solving $h = (1/2)at^2$ to yield

$$\begin{aligned} t &= \sqrt{2h/a} \\ &= w \sqrt{\frac{\mu}{2T}}. \end{aligned}$$



n / A triangular pulse spreads out.

Our final result for the speed of the pulses is

$$\begin{aligned} v &= w/t \\ &= \sqrt{\frac{2T}{\mu}}. \end{aligned}$$

The remarkable feature of this result is that the velocity of the pulses does not depend at all on w or h , i.e., any triangular pulse has the same speed. It is an experimental fact (and we will also prove rigorously below) that any pulse of any kind, triangular or otherwise, travels along the string at the same speed. Of course, after so many approximations we cannot expect to have gotten all the numerical factors right. The correct result for the speed of the pulses is

$$v = \sqrt{\frac{T}{\mu}}.$$

The importance of the above derivation lies in the insight it brings — that all pulses move with the same speed — rather than in the details of the numerical result. The reason for our too-high value for the velocity is not hard to guess. It comes from the assumption that the acceleration was constant, when actually the total force on the segment would diminish as it flattened out.

Treatment using calculus

After expending considerable effort for an approximate solution, we now display the power of calculus with a rigorous and completely general treatment that is nevertheless much shorter and easier. Let the flat position of the string define the x axis, so that y measures how far a point on the string is from equilibrium. The motion of the string is characterized by $y(x, t)$, a function of two variables. Knowing that the force on any small segment of string depends on the curvature of the string in that area, and that the second derivative is a measure of curvature, it is not surprising to find that the infinitesimal force dF acting on an infinitesimal segment dx is given by

$$dF = T \frac{\partial^2 y}{\partial x^2} dx.$$

(This can be proved by vector addition of the two infinitesimal forces acting on either side.) The symbol ∂ stands for a partial derivative, e.g., $\partial/\partial x$ means a derivative with respect to x that is evaluated while treating t as a constant. The acceleration is then $a = dF/dm$, or, substituting $dm = \mu dx$,

$$\frac{\partial^2 y}{\partial t^2} = \frac{T}{\mu} \frac{\partial^2 y}{\partial x^2}.$$

The second derivative with respect to time is related to the second derivative with respect to position. This is no more than a fancy

mathematical statement of the intuitive fact developed above, that the string accelerates so as to flatten out its curves.

Before even bothering to look for solutions to this equation, we note that it already proves the principle of superposition, because the derivative of a sum is the sum of the derivatives. Therefore the sum of any two solutions will also be a solution.

Based on experiment, we expect that this equation will be satisfied by any function $y(x, t)$ that describes a pulse or wave pattern moving to the left or right at the correct speed v . In general, such a function will be of the form $y = f(x - vt)$ or $y = f(x + vt)$, where f is any function of one variable. Because of the chain rule, each derivative with respect to time brings out a factor of v . Evaluating the second derivatives on both sides of the equation gives

$$(\pm v)^2 f'' = \frac{T}{\mu} f''.$$

Squaring gets rid of the sign, and we find that we have a valid solution for any function f , provided that v is given by

$$v = \sqrt{\frac{T}{\mu}}.$$

6.1.3 Sound and light waves

Sound waves

The phenomenon of sound is easily found to have all the characteristics we expect from a wave phenomenon:

- Sound waves obey superposition. Sounds do not knock other sounds out of the way when they collide, and we can hear more than one sound at once if they both reach our ear simultaneously.
- The medium does not move with the sound. Even standing in front of a titanic speaker playing earsplitting music, we do not feel the slightest breeze.
- The velocity of sound depends on the medium. Sound travels faster in helium than in air, and faster in water than in helium. Putting more energy into the wave makes it more intense, not faster. For example, you can easily detect an echo when you clap your hands a short distance from a large, flat wall, and the delay of the echo is no shorter for a louder clap.

Although not all waves have a speed that is independent of the shape of the wave, and this property therefore is irrelevant to our collection of evidence that sound is a wave phenomenon, sound does nevertheless have this property. For instance, the music in a large concert hall or stadium may take on the order of a second to reach someone seated in the nosebleed section, but we do not notice or

care, because the delay is the same for every sound. Bass, drums, and vocals all head outward from the stage at 340 m/s, regardless of their differing wave shapes. (The speed of sound in a gas is related to the gas's physical properties in example 13 on p. 389.)

If sound has all the properties we expect from a wave, then what type of wave is it? It is a series of compressions and expansions of the air. Even for a very loud sound, the increase or decrease compared to normal atmospheric pressure is no more than a part per million, so our ears are apparently very sensitive instruments. In a vacuum, there is no medium for the sound waves, and so they cannot exist. The roars and whooshes of space ships in Hollywood movies are fun, but scientifically wrong.

Light waves

Entirely similar observations lead us to believe that light is a wave, although the concept of light as a wave had a long and tortuous history. It is interesting to note that Isaac Newton very influentially advocated a contrary idea about light. The belief that matter was made of atoms was stylish at the time among radical thinkers (although there was no experimental evidence for their existence), and it seemed logical to Newton that light as well should be made of tiny particles, which he called corpuscles (Latin for "small objects"). Newton's triumphs in the science of mechanics, i.e., the study of matter, brought him such great prestige that nobody bothered to question his incorrect theory of light for 150 years. One persuasive proof that light is a wave is that according to Newton's theory, two intersecting beams of light should experience at least some disruption because of collisions between their corpuscles. Even if the corpuscles were extremely small, and collisions therefore very infrequent, at least some dimming should have been measurable. In fact, very delicate experiments have shown that there is no dimming.

The wave theory of light was entirely successful up until the 20th century, when it was discovered that not all the phenomena of light could be explained with a pure wave theory. It is now believed that both light and matter are made out of tiny chunks which have both wave and particle properties. For now, we will content ourselves with the wave theory of light, which is capable of explaining a great many things, from cameras to rainbows.

If light is a wave, what is waving? What is the medium that wiggles when a light wave goes by? It isn't air. A vacuum is impenetrable to sound, but light from the stars travels happily through billions of miles of empty space. Light bulbs have no air inside them, but that doesn't prevent the light waves from leaving the filament. For a long time, physicists assumed that there must be a mysterious medium for light waves, and they called it the ether (not to be confused with the chemical). Supposedly the ether existed everywhere in space, and was immune to vacuum pumps. The details of

the story are more fittingly reserved for later in this course, but the end result was that a long series of experiments failed to detect any evidence for the ether, and it is no longer believed to exist. Instead, light can be explained as a wave pattern made up of electrical and magnetic fields.

6.1.4 Periodic waves

Period and frequency of a periodic wave

You choose a radio station by selecting a certain frequency. We have already defined period and frequency for vibrations,

$$T = \text{period} = \text{seconds per cycle}$$

$$f = \text{frequency} = 1/T = \text{cycles per second}$$

$$\omega = \text{angular frequency} = 2\pi f = \text{radians per second}$$

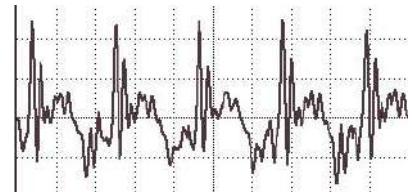
but what do they signify in the case of a wave? We can recycle our previous definition simply by stating it in terms of the vibrations that the wave causes as it passes a receiving instrument at a certain point in space. For a sound wave, this receiver could be an eardrum or a microphone. If the vibrations of the eardrum repeat themselves over and over, i.e., are periodic, then we describe the sound wave that caused them as periodic. Likewise we can define the period and frequency of a wave in terms of the period and frequency of the vibrations it causes. As another example, a periodic water wave would be one that caused a rubber duck to bob in a periodic manner as they passed by it.

The period of a sound wave correlates with our sensory impression of musical pitch. A high frequency (short period) is a high note. The sounds that really define the musical notes of a song are only the ones that are periodic. It is not possible to sing a nonperiodic sound like “sh” with a definite pitch.

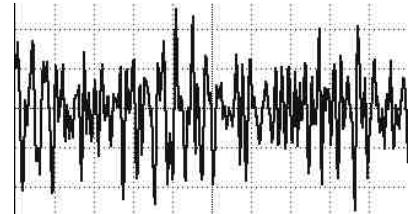
The frequency of a light wave corresponds to color. Violet is the high-frequency end of the rainbow, red the low-frequency end. A color like brown that does not occur in a rainbow is not a periodic light wave. Many phenomena that we do not normally think of as light are actually just forms of light that are invisible because they fall outside the range of frequencies our eyes can detect. Beyond the red end of the visible rainbow, there are infrared and radio waves. Past the violet end, we have ultraviolet, x-rays, and gamma rays.

Graphs of waves as a function of position

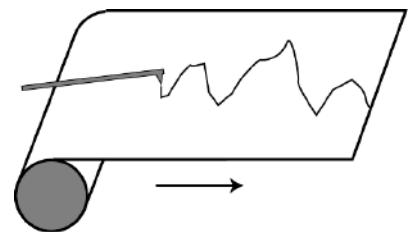
Some waves, like sound waves, are easy to study by placing a detector at a certain location in space and studying the motion as a function of time. The result is a graph whose horizontal axis is time. With a water wave, on the other hand, it is simpler just to look at the wave directly. This visual snapshot amounts to a graph of the height of the water wave as a function of position. Any wave can be represented in either way.



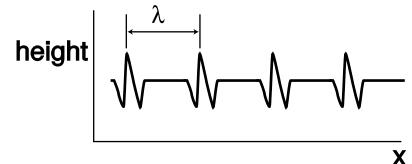
o / A graph of pressure versus time for a periodic sound wave, the vowel “ah.”



p / A similar graph for a non-periodic wave, “sh.”



q / A strip chart recorder.



r / A water wave profile created by a series of repeating pulses.

An easy way to visualize this is in terms of a strip chart recorder, an obsolescing device consisting of a pen that wiggles back and forth as a roll of paper is fed under it. It can be used to record a person's electrocardiogram, or seismic waves too small to be felt as a noticeable earthquake but detectable by a seismometer. Taking the seismometer as an example, the chart is essentially a record of the ground's wave motion as a function of time, but if the paper was set to feed at the same velocity as the motion of an earthquake wave, it would also be a full-scale representation of the profile of the actual wave pattern itself. Assuming, as is usually the case, that the wave velocity is a constant number regardless of the wave's shape, knowing the wave motion as a function of time is equivalent to knowing it as a function of position.

Wavelength

Any wave that is periodic will also display a repeating pattern when graphed as a function of position. The distance spanned by one repetition is referred to as one wavelength. The usual notation for wavelength is λ , the Greek letter lambda. Wavelength is to space as period is to time.

Wave velocity related to frequency and wavelength

Suppose that we create a repetitive disturbance by kicking the surface of a swimming pool. We are essentially making a series of wave pulses. The wavelength is simply the distance a pulse is able to travel before we make the next pulse. The distance between pulses is λ , and the time between pulses is the period, T , so the speed of the wave is the distance divided by the time,

$$v = \lambda/T.$$

This important and useful relationship is more commonly written in terms of the frequency,

$$v = f\lambda.$$

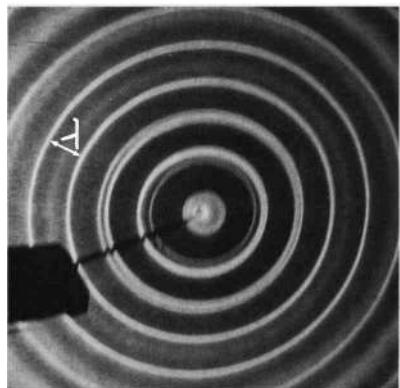
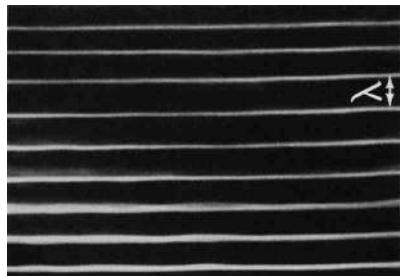
Wavelength of radio waves

example 5

- ▷ The speed of light is 3.0×10^8 m/s. What is the wavelength of the radio waves emitted by KMHD, a station whose frequency is 89.1 MHz?

▷ Solving for wavelength, we have

$$\begin{aligned} \lambda &= v/f \\ &= (3.0 \times 10^8 \text{ m/s}) / (89.1 \times 10^6 \text{ s}^{-1}) \\ &= 3.4 \text{ m} \end{aligned}$$



s / Wavelengths of linear and circular waves.

The size of a radio antenna is closely related to the wavelength of the waves it is intended to receive. The match need not be exact

(since after all one antenna can receive more than one wavelength!), but the ordinary “whip” antenna such as a car’s is $1/4$ of a wavelength. An antenna optimized to receive KMHD’s signal would have a length of $(3.4 \text{ m})/4 = 0.85 \text{ m}$.

The equation $v = f\lambda$ defines a fixed relationship between any two of the variables if the other is held fixed. The speed of radio waves in air is almost exactly the same for all wavelengths and frequencies (it is exactly the same if they are in a vacuum), so there is a fixed relationship between their frequency and wavelength. Thus we can say either “Are we on the same wavelength?” or “Are we on the same frequency?”

A different example is the behavior of a wave that travels from a region where the medium has one set of properties to an area where the medium behaves differently. The frequency is now fixed, because otherwise the two portions of the wave would otherwise get out of step, causing a kink or discontinuity at the boundary, which would be unphysical. (A more careful argument is that a kink or discontinuity would have infinite curvature, and waves tend to flatten out their curvature. An infinite curvature would flatten out infinitely fast, i.e., it could never occur in the first place.) Since the frequency must stay the same, any change in the velocity that results from the new medium must cause a change in wavelength.

The velocity of water waves depends on the depth of the water, so based on $\lambda = v/f$, we see that water waves that move into a region of different depth must change their wavelength, as shown in figure u. This effect can be observed when ocean waves come up to the shore. If the deceleration of the wave pattern is sudden enough, the tip of the wave can curl over, resulting in a breaking wave.

A note on dispersive waves

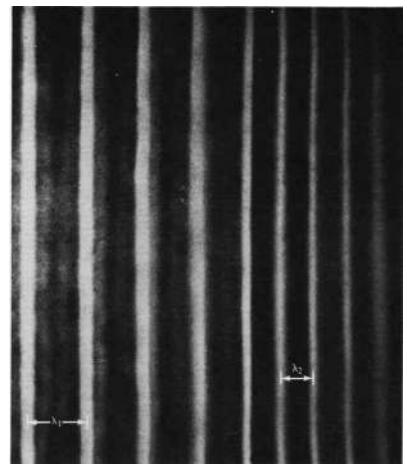
The discussion of wave velocity given here is actually a little bit of an oversimplification for a wave whose velocity depends on its frequency and wavelength. Such a wave is called a dispersive wave. Nearly all the waves we deal with in this course are nondispersive, but the issue becomes important in chapter 13, where it is discussed in detail.

Sinusoidal waves

Sinusoidal waves are the most important special case of periodic waves. In fact, many scientists and engineers would be uncomfortable with defining a waveform like the “ah” vowel sound as having a definite frequency and wavelength, because they consider only sine waves to be pure examples of a certain frequency and wavelengths. Their bias is not unreasonable, since the French mathematician Fourier showed that any periodic wave with frequency f can be constructed as a superposition of sine waves with frequencies $f, 2f, 3f, \dots$. In this sense, sine waves are the basic, pure building



t / Ultrasound, i.e., sound with frequencies higher than the range of human hearing, was used to make this image of a fetus. The resolution of the image is related to the wavelength, since details smaller than about one wavelength cannot be resolved. High resolution therefore requires a short wavelength, corresponding to a high frequency.



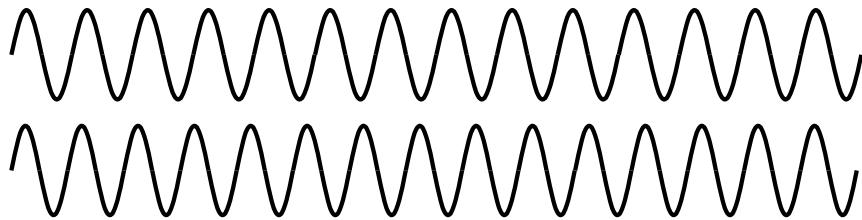
u / A water wave traveling into a region with different depth will change its wavelength.

blocks of all waves. (Fourier's result so surprised the mathematical community of France that he was ridiculed the first time he publicly presented his theorem.)

However, what definition to use is really a matter of convenience. Our sense of hearing perceives any two sounds having the same period as possessing the same pitch, regardless of whether they are sine waves or not. This is undoubtedly because our ear-brain system evolved to be able to interpret human speech and animal noises, which are periodic but not sinusoidal. Our eyes, on the other hand, judge a color as pure (belonging to the rainbow set of colors) only if it is a sine wave.

Discussion Questions

A Suppose we superimpose two sine waves with equal amplitudes but slightly different frequencies, as shown in the figure. What will the superposition look like? What would this sound like if they were sound waves?



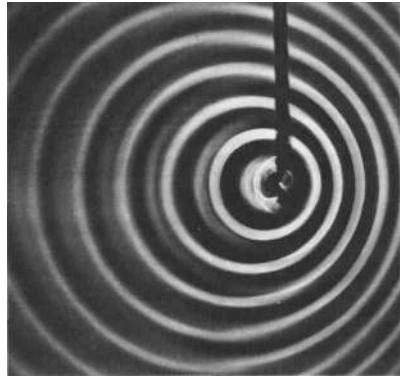
v / Discussion question A.

6.1.5 The Doppler effect

Figure w shows the wave pattern made by the tip of a vibrating rod which is moving across the water. If the rod had been vibrating in one place, we would have seen the familiar pattern of concentric circles, all centered on the same point. But since the source of the waves is moving, the wavelength is shortened on one side and lengthened on the other. This is known as the Doppler effect.

Note that the velocity of the waves is a fixed property of the medium, so for example the forward-going waves do not get an extra boost in speed as would a material object like a bullet being shot forward from an airplane.

We can also infer a change in frequency. Since the velocity is constant, the equation $v = f\lambda$ tells us that the change in wavelength must be matched by an opposite change in frequency: higher frequency for the waves emitted forward, and lower for the ones emitted backward. The frequency Doppler effect is the reason for the familiar dropping-pitch sound of a race car going by. As the car approaches us, we hear a higher pitch, but after it passes us we hear



w / The pattern of waves made by a point source moving to the right across the water. Note the shorter wavelength of the forward-emitted waves and the longer wavelength of the backward-going ones.

a frequency that is lower than normal.

The Doppler effect will also occur if the observer is moving but the source is stationary. For instance, an observer moving toward a stationary source will perceive one crest of the wave, and will then be surrounded by the next crest sooner than she otherwise would have, because she has moved toward it and hastened her encounter with it. Roughly speaking, the Doppler effect depends only the relative motion of the source and the observer, not on their absolute state of motion (which is not a well-defined notion in physics) or on their velocity relative to the medium.

Restricting ourselves to the case of a moving source, and to waves emitted either directly along or directly against the direction of motion, we can easily calculate the wavelength, or equivalently the frequency, of the Doppler-shifted waves. Let u be the velocity of the source. The wavelength of the forward-emitted waves is shortened by an amount uT equal to the distance traveled by the source over the course of one period. Using the definition $f = 1/T$ and the equation $v = f\lambda$, we find for the wavelength λ' of the Doppler-shifted wave the equation

$$\lambda' = \left(1 - \frac{u}{v}\right) \lambda.$$

A similar equation can be used for the backward-emitted waves, but with a plus sign rather than a minus sign.

Doppler-shifted sound from a race car *example 6*

▷ If a race car moves at a velocity of 50 m/s, and the velocity of sound is 340 m/s, by what percentage are the wavelength and frequency of its sound waves shifted for an observer lying along its line of motion?

▷ For an observer whom the car is approaching, we find

$$1 - \frac{u}{v} = 0.85,$$

so the shift in wavelength is 15%. Since the frequency is inversely proportional to the wavelength for a fixed value of the speed of sound, the frequency is shifted upward by

$$1/0.85 = 1.18,$$

i.e., a change of 18%. (For velocities that are small compared to the wave velocities, the Doppler shifts of the wavelength and frequency are about the same.)

Doppler shift of the light emitted by a race car *example 7*

▷ What is the percent shift in the wavelength of the light waves emitted by a race car's headlights?

▷ Looking up the speed of light in the back of the book, $v = 3.0 \times 10^8$ m/s, we find

$$1 - \frac{u}{v} = 0.99999983,$$



x / The galaxy M100 in the constellation Coma Berenices. Under higher magnification, the milky clouds reveal themselves to be composed of trillions of stars.



y / The telescope at Mount Wilson used by Hubble.

i.e., the percentage shift is only 0.000017%.

The second example shows that under ordinary earthbound circumstances, Doppler shifts of light are negligible because ordinary things go so much slower than the speed of light. It's a different story, however, when it comes to stars and galaxies, and this leads us to a story that has profound implications for our understanding of the origin of the universe.

The Big Bang

As soon as astronomers began looking at the sky through telescopes, they began noticing certain objects that looked like clouds in deep space. The fact that they looked the same night after night meant that they were beyond the earth's atmosphere. Not knowing what they really were, but wanting to sound official, they called them "nebulae," a Latin word meaning "clouds" but sounding more impressive. In the early 20th century, astronomers realized that although some really were clouds of gas (e.g., the middle "star" of Orion's sword, which is visibly fuzzy even to the naked eye when conditions are good), others were what we now call galaxies: virtual island universes consisting of trillions of stars (for example the Andromeda Galaxy, which is visible as a fuzzy patch through binoculars). Three hundred years after Galileo had resolved the Milky Way into individual stars through his telescope, astronomers realized that the universe is made of galaxies of stars, and the Milky Way is simply the visible part of the flat disk of our own galaxy, seen from inside.

This opened up the scientific study of cosmology, the structure and history of the universe as a whole, a field that had not been seriously attacked since the days of Newton. Newton had realized that if gravity was always attractive, never repulsive, the universe would have a tendency to collapse. His solution to the problem was to posit a universe that was infinite and uniformly populated with matter, so that it would have no geometrical center. The gravitational forces in such a universe would always tend to cancel out by symmetry, so there would be no collapse. By the 20th century, the belief in an unchanging and infinite universe had become conventional wisdom in science, partly as a reaction against the time that had been wasted trying to find explanations of ancient geological phenomena based on catastrophes suggested by biblical events like Noah's flood.

In the 1920's astronomer Edwin Hubble began studying the Doppler shifts of the light emitted by galaxies. A former college football player with a serious nicotine addiction, Hubble did not set out to change our image of the beginning of the universe. His autobiography seldom even mentions the cosmological discovery for which he is now remembered. When astronomers began to study the



Doppler shifts of galaxies, they expected that each galaxy's direction and velocity of motion would be essentially random. Some would be approaching us, and their light would therefore be Doppler-shifted to the blue end of the spectrum, while an equal number would be expected to have red shifts. What Hubble discovered instead was that except for a few very nearby ones, all the galaxies had red shifts, indicating that they were receding from us at a hefty fraction of the speed of light. Not only that, but the ones farther away were receding more quickly. The speeds were directly proportional to their distance from us.

Did this mean that the earth (or at least our galaxy) was the center of the universe? No, because Doppler shifts of light only depend on the relative motion of the source and the observer. If we see a distant galaxy moving away from us at 10% of the speed of light, we can be assured that the astronomers who live in that galaxy will see ours receding from them at the same speed in the opposite direction. The whole universe can be envisioned as a rising loaf of raisin bread. As the bread expands, there is more and more space between the raisins. The farther apart two raisins are, the greater the speed with which they move apart.

The universe's expansion is presumably decelerating because of gravitational attraction among the galaxies. We do not presently know whether there is enough mass in the universe to cause enough attraction to halt the expansion eventually. But perhaps more interesting than the distant future of the universe is what its present expansion implies about its past. Extrapolating backward in time using the known laws of physics, the universe must have been denser and denser at earlier and earlier times. At some point, it must have been extremely dense and hot, and we can even detect the radiation from this early fireball, in the form of microwave radiation that permeates space. The phrase Big Bang was originally coined by the doubters of the theory to make it sound ridiculous, but it stuck, and today essentially all astronomers accept the Big Bang theory based on the very direct evidence of the red shifts and the cosmic microwave background radiation.

Finally it should be noted what the Big Bang theory is not. It is not an explanation of *why* the universe exists. Such questions belong to the realm of religion, not science. Science can find ever simpler and ever more fundamental explanations for a variety of phenomena, but ultimately science takes the universe as it is according to observations.

Furthermore, there is an unfortunate tendency, even among many scientists, to speak of the Big Bang theory as a description of the very first event in the universe, which caused everything after it. Although it is true that time may have had a beginning (Einstein's theory of general relativity admits such a possibility), the methods

z / How do astronomers know what mixture of wavelengths a star emitted originally, so that they can tell how much the Doppler shift was? This image (obtained by the author with equipment costing about \$5, and no telescope) shows the mixture of colors emitted by the star Sirius. (If you have the book in black and white, blue is on the left and red on the right.) The star appears white or bluish-white to the eye, but any light looks white if it contains roughly an equal mixture of the rainbow colors, i.e., of all the pure sinusoidal waves with wavelengths lying in the visible range. Note the black "gap teeth." These are the fingerprint of hydrogen in the outer atmosphere of Sirius. These wavelengths are selectively absorbed by hydrogen. Sirius is in our own galaxy, but similar stars in other galaxies would have the whole pattern shifted toward the red end, indicating they are moving away from us.

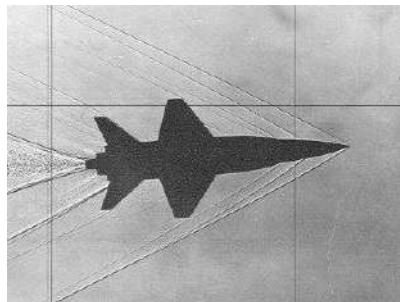
of science can only work within a certain range of conditions such as temperature and density. Beyond a temperature of about 10^9 K, the random thermal motion of subatomic particles becomes so rapid that its velocity is comparable to the speed of light. Early enough in the history of the universe, when these temperatures existed, Newtonian physics becomes less accurate, and we must describe nature using the more general description given by Einstein's theory of relativity, which encompasses Newtonian physics as a special case. At even higher temperatures, beyond about 10^{33} degrees, physicists know that Einstein's theory as well begins to fall apart, but we don't know how to construct the even more general theory of nature that would work at those temperatures. No matter how far physics progresses, we will never be able to describe nature at infinitely high temperatures, since there is a limit to the temperatures we can explore by experiment and observation in order to guide us to the right theory. We are confident that we understand the basic physics involved in the evolution of the universe starting a few minutes after the Big Bang, and we may be able to push back to milliseconds or microseconds after it, but we cannot use the methods of science to deal with the beginning of time itself.

A note on Doppler shifts of light

If Doppler shifts depend only on the relative motion of the source and receiver, then there is no way for a person moving with the source and another person moving with the receiver to determine who is moving and who isn't. Either can blame the Doppler shift entirely on the other's motion and claim to be at rest herself. This is entirely in agreement with the principle stated originally by Galileo that all motion is relative.

On the other hand, a careful analysis of the Doppler shifts of water or sound waves shows that it is only approximately true, at low speeds, that the shifts just depend on the relative motion of the source and observer. For instance, it is possible for a jet plane to keep up with its own sound waves, so that the sound waves appear to stand still to the pilot of the plane. The pilot then knows she is moving at exactly the speed of sound. The reason this doesn't disprove the relativity of motion is that the pilot is not really determining her absolute motion but rather her motion relative to the air, which is the medium of the sound waves.

Einstein realized that this solved the problem for sound or water waves, but would not salvage the principle of relative motion in the case of light waves, since light is not a vibration of any physical medium such as water or air. Beginning by imagining what a beam of light would look like to a person riding a motorcycle alongside it, Einstein eventually came up with a radical new way of describing the universe, in which space and time are distorted as measured by observers in different states of motion. As a consequence of this



aa / Shock waves are created by the X-15 rocket plane, flying at 3.5 times the speed of sound.



ab / This fighter jet has just accelerated past the speed of sound. The sudden decompression of the air causes water droplets to condense, forming a cloud.

Theory of Relativity, he showed that light waves would have Doppler shifts that would exactly, not just approximately, depend only on the relative motion of the source and receiver.

Discussion Questions

A If an airplane travels at exactly the speed of sound, what would be the wavelength of the forward-emitted part of the sound waves it emitted? How should this be interpreted, and what would actually happen? What happens if it's going faster than the speed of sound? Can you use this to explain what you see in figures aa and ab?

B If bullets go slower than the speed of sound, why can a supersonic fighter plane catch up to its own sound, but not to its own bullets?

C If someone inside a plane is talking to you, should their speech be Doppler shifted?

6.2 Bounded waves

Speech is what separates humans most decisively from animals. No other species can master syntax, and even though chimpanzees can learn a vocabulary of hand signs, there is an unmistakable difference between a human infant and a baby chimp: starting from birth, the human experiments with the production of complex speech sounds.

Since speech sounds are instinctive for us, we seldom think about them consciously. How do we control sound waves so skillfully? Mostly we do it by changing the shape of a connected set of hollow cavities in our chest, throat, and head. Somehow by moving the boundaries of this space in and out, we can produce all the vowel sounds. Up until now, we have been studying only those properties of waves that can be understood as if they existed in an infinite, open space with no boundaries. In this chapter we address what happens when a wave is confined within a certain space, or when a wave pattern encounters the boundary between two different media, such as when a light wave moving through air encounters a glass windowpane.

6.2.1 Reflection, transmission, and absorption

Reflection and transmission

Sound waves can echo back from a cliff, and light waves are reflected from the surface of a pond. We use the word reflection, normally applied only to light waves in ordinary speech, to describe any such case of a wave rebounding from a barrier. Figure (a) shows a circular water wave being reflected from a straight wall. In this chapter, we will concentrate mainly on reflection of waves that move in one dimension, as in figure c/1.

Wave reflection does not surprise us. After all, a material object such as a rubber ball would bounce back in the same way. But waves are not objects, and there are some surprises in store.

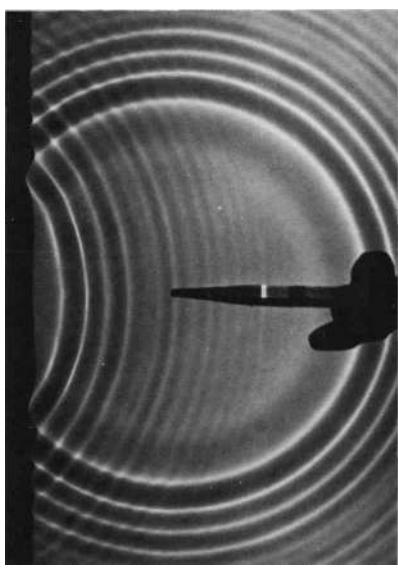
First, only part of the wave is usually reflected. Looking out through a window, we see light waves that passed through it, but a person standing outside would also be able to see her reflection in the glass. A light wave that strikes the glass is partly reflected and partly transmitted (passed) by the glass. The energy of the original wave is split between the two. This is different from the behavior of the rubber ball, which must go one way or the other, not both.

Second, consider what you see if you are swimming underwater and you look up at the surface. You see your own reflection. This is utterly counterintuitive, since we would expect the light waves to burst forth to freedom in the wide-open air. A material projectile shot up toward the surface would never rebound from the water-air boundary!

What is it about the difference between two media that causes

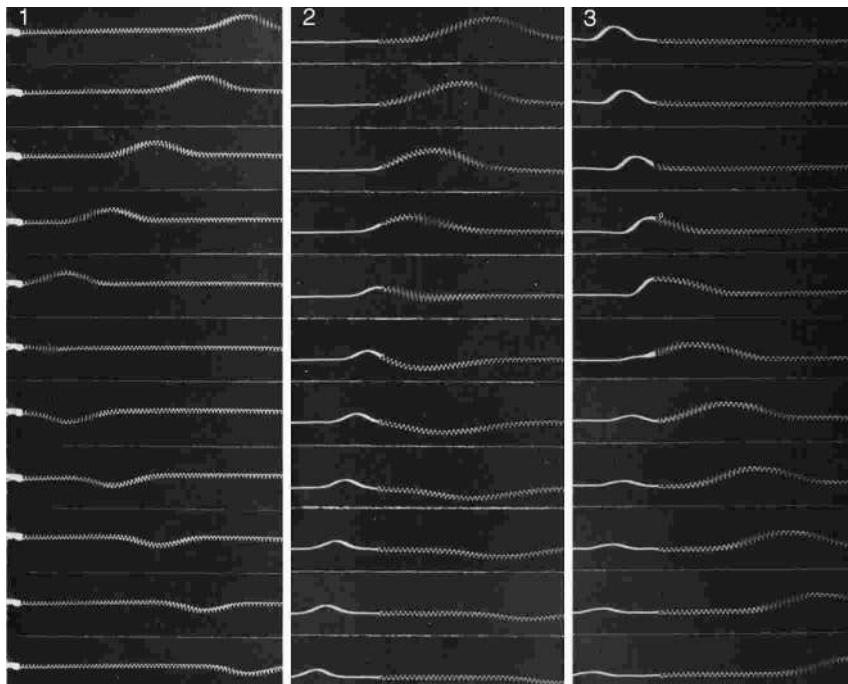


a / A cross-sectional view of a human body, showing the vocal tract.



b / Circular water waves are reflected from a boundary on the left. (*PSSC Physics*)

waves to be partly reflected at the boundary between them? Is it their density? Their chemical composition? Typically all that matters is the speed of the wave in the two media.² A wave is partially reflected and partially transmitted at the boundary between media in which it has different speeds. For example, the speed of light waves in window glass is about 30% less than in air, which explains why windows always make reflections. Figure c shows examples of wave pulses being reflected at the boundary between two coil springs of different weights, in which the wave speed is different.



c / 1. A wave on a coil spring, initially traveling to the left, is reflected from the fixed end. 2. A wave in the lighter spring, where the wave speed is greater, travels to the left and is then partly reflected and partly transmitted at the boundary with the heavier coil spring, which has a lower wave speed. The reflection is inverted. 3. A wave moving to the right in the heavier spring is partly reflected at the boundary with the lighter spring. The reflection is uninverted. (PSSC Physics)

Reflections such as b and c/1, where a wave encounters a massive fixed object, can usually be understood on the same basis as cases like c/2 and c/3 where two media meet. Example c/1, for instance, is like a more extreme version of example c/2. If the heavy coil spring in c/2 was made heavier and heavier, it would end up acting like the fixed wall to which the light spring in c/1 has been attached.

self-check B

In figure c/1, the reflected pulse is upside-down, but its depth is just as big as the original pulse's height. How does the energy of the reflected pulse compare with that of the original? ▷ Answer, p. 1061

Fish have internal ears.

example 8

Why don't fish have ear-holes? The speed of sound waves in a fish's body is not much different from their speed in water, so sound waves are not strongly reflected from a fish's skin. They pass right through its body, so fish can have internal ears.

²Some exceptions are described in sec. 6.2.5, p. 389.

Whale songs traveling long distances

example 9

Sound waves travel at drastically different speeds through rock, water, and air. Whale songs are thus strongly reflected both at both the bottom and the surface. The sound waves can travel hundreds of miles, bouncing repeatedly between the bottom and the surface, and still be detectable. Sadly, noise pollution from ships has nearly shut down this cetacean version of the internet.

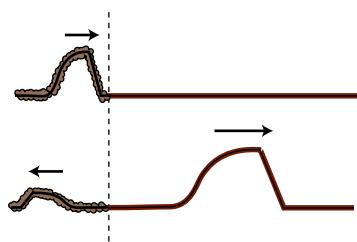
Long-distance radio communication

example 10

Radio communication can occur between stations on opposite sides of the planet. The mechanism is entirely similar to the one explained in the previous example, but the three media involved are the earth, the atmosphere, and the ionosphere.

self-check C

Sonar is a method for ships and submarines to detect each other by producing sound waves and listening for echoes. What properties would an underwater object have to have in order to be invisible to sonar? ▷ Answer, p. 1061

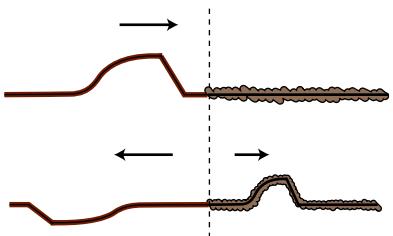


d / An uninverted reflection. The reflected pulse is reversed front to back, but is not upside-down.

The use of the word “reflection” naturally brings to mind the creation of an image by a mirror, but this might be confusing, because we do not normally refer to “reflection” when we look at surfaces that are not shiny. Nevertheless, reflection is how we see the surfaces of all objects, not just polished ones. When we look at a sidewalk, for example, we are actually seeing the reflecting of the sun from the concrete. The reason we don’t see an image of the sun at our feet is simply that the rough surface blurs the image so drastically.

Inverted and uninverted reflections

Notice how the pulse reflected back to the right in example c/2 comes back upside-down, whereas the one reflected back to the left in c/3 returns in its original upright form. This is true for other waves as well. In general, there are two possible types of reflections, a reflection back into a faster medium and a reflection back into a slower medium. One type will always be an inverting reflection and one noninverting.



e / An inverted reflection. The reflected pulse is reversed both front to back and top to bottom.

It’s important to realize that when we discuss inverted and uninverted reflections on a string, we are talking about whether the wave is flipped across the direction of motion (i.e., upside-down in these drawings). The reflected pulse will always be reversed front to back, as shown in figures d and e. This is because it is traveling in the other direction. The leading edge of the pulse is what gets reflected first, so it is still ahead when it starts back to the left — it’s just that “ahead” is now in the opposite direction.

Absorption

So far we have tacitly assumed that wave energy remains as wave energy, and is not converted to any other form. If this was true, then the world would become more and more full of sound waves, which could never escape into the vacuum of outer space. In reality, any mechanical wave consists of a traveling pattern of vibrations of some physical medium, and vibrations of matter always produce heat, as when you bend a coathanger back and forth and it becomes hot. We can thus expect that in mechanical waves such as water waves, sound waves, or waves on a string, the wave energy will gradually be converted into heat. This is referred to as absorption. The reduction in the wave's energy can also be described as a reduction in amplitude, the relationship between them being, as with a vibrating object, $E \propto A^2$.

The wave suffers a decrease in amplitude, as shown in figure f. The decrease in amplitude amounts to the same fractional change for each unit of distance covered. For example, if a wave decreases from amplitude 2 to amplitude 1 over a distance of 1 meter, then after traveling another meter it will have an amplitude of 1/2. That is, the reduction in amplitude is exponential. This can be proved as follows. By the principle of superposition, we know that a wave of amplitude 2 must behave like the superposition of two identical waves of amplitude 1. If a single amplitude-1 wave would die down to amplitude 1/2 over a certain distance, then two amplitude-1 waves superposed on top of one another to make amplitude 1+1=2 must die down to amplitude 1/2+1/2=1 over the same distance.

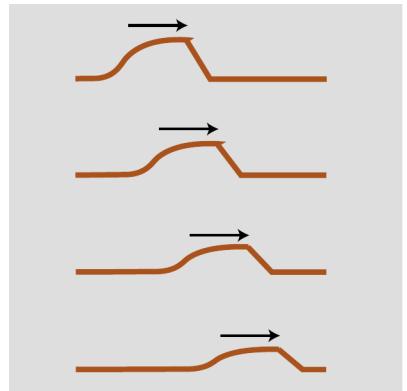
self-check D

As a wave undergoes absorption, it loses energy. Does this mean that it slows down?

▷ Answer, p. 1061

In many cases, this frictional heating effect is quite weak. Sound waves in air, for instance, dissipate into heat extremely slowly, and the sound of church music in a cathedral may reverberate for as much as 3 or 4 seconds before it becomes inaudible. During this time it has traveled over a kilometer! Even this very gradual dissipation of energy occurs mostly as heating of the church's walls and by the leaking of sound to the outside (where it will eventually end up as heat). Under the right conditions (humid air and low frequency), a sound wave in a straight pipe could theoretically travel hundreds of kilometers before being noticeable attenuated.

In general, the absorption of mechanical waves depends a great deal on the chemical composition and microscopic structure of the medium. Ripples on the surface of antifreeze, for instance, die out extremely rapidly compared to ripples on water. For sound waves and surface waves in liquids and gases, what matters is the viscosity of the substance, i.e., whether it flows easily like water or mercury or more sluggishly like molasses or antifreeze. This explains why



f / A pulse traveling through a highly absorptive medium.

our intuitive expectation of strong absorption of sound in water is incorrect. Water is a very weak absorber of sound (viz. whale songs and sonar), and our incorrect intuition arises from focusing on the wrong property of the substance: water's high density, which is irrelevant, rather than its low viscosity, which is what matters.

Light is an interesting case, since although it can travel through matter, it is not itself a vibration of any material substance. Thus we can look at the star Sirius, 10^{14} km away from us, and be assured that none of its light was absorbed in the vacuum of outer space during its 9-year journey to us. The Hubble Space Telescope routinely observes light that has been on its way to us since the early history of the universe, billions of years ago. Of course the energy of light can be dissipated if it does pass through matter (and the light from distant galaxies is often absorbed if there happen to be clouds of gas or dust in between).

Soundproofing

example 11

Typical amateur musicians setting out to soundproof their garages tend to think that they should simply cover the walls with the densest possible substance. In fact, sound is not absorbed very strongly even by passing through several inches of wood. A better strategy for soundproofing is to create a sandwich of alternating layers of materials in which the speed of sound is very different, to encourage reflection.

The classic design is alternating layers of fiberglass and plywood. The speed of sound in plywood is very high, due to its stiffness, while its speed in fiberglass is essentially the same as its speed in air. Both materials are fairly good sound absorbers, but sound waves passing through a few inches of them are still not going to be absorbed sufficiently. The point of combining them is that a sound wave that tries to get out will be strongly reflected at each of the fiberglass-plywood boundaries, and will bounce back and forth many times like a ping pong ball. Due to all the back-and-forth motion, the sound may end up traveling a total distance equal to ten times the actual thickness of the soundproofing before it escapes. This is the equivalent of having ten times the thickness of sound-absorbing material.

Radio transmission

example 12

A radio transmitting station must have a length of wire or cable connecting the amplifier to the antenna. The cable and the antenna act as two different media for radio waves, and there will therefore be partial reflection of the waves as they come from the cable to the antenna. If the waves bounce back and forth many times between the amplifier and the antenna, a great deal of their energy will be absorbed. There are two ways to attack the problem. One possibility is to design the antenna so that the speed of the waves in it is as close as possible to the speed of the waves

in the cable; this minimizes the amount of reflection. The other method is to connect the amplifier to the antenna using a type of wire or cable that does not strongly absorb the waves. Partial reflection then becomes irrelevant, since all the wave energy will eventually exit through the antenna.

Discussion Questions

A A sound wave that underwent a pressure-inverting reflection would have its compressions converted to expansions and vice versa. How would its energy and frequency compare with those of the original sound? Would it sound any different? What happens if you swap the two wires where they connect to a stereo speaker, resulting in waves that vibrate in the opposite way?

6.2.2 Quantitative treatment of reflection

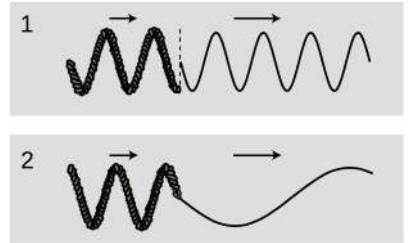
In this section we use the example of waves on a string to analyze the reasons why a reflection occurs at the boundary between media, predict quantitatively the intensities of reflection and transmission, and discuss how to tell which reflections are inverting and which are noninverting. Some technical details are relegated to sec. 6.2.5, p. 389.

Why reflection occurs

To understand the fundamental reasons for what does occur at the boundary between media, let's first discuss what doesn't happen. For the sake of concreteness, consider a sinusoidal wave on a string. If the wave progresses from a heavier portion of the string, in which its velocity is low, to a lighter-weight part, in which it is high, then the equation $v = f\lambda$ tells us that it must change its frequency, or its wavelength, or both. If only the frequency changed, then the parts of the wave in the two different portions of the string would quickly get out of step with each other, producing a discontinuity in the wave, g/1. This is unphysical, so we know that the wavelength must change while the frequency remains constant, g/2.

But there is still something unphysical about figure g/2. The sudden change in the shape of the wave has resulted in a sharp kink at the boundary. This can't really happen, because the medium tends to accelerate in such a way as to eliminate curvature. A sharp kink corresponds to an infinite curvature at one point, which would produce an infinite acceleration, which would not be consistent with the smooth pattern of wave motion envisioned in fig. g/2. Waves can have kinks, but not stationary kinks.

We conclude that without positing partial reflection of the wave, we cannot simultaneously satisfy the requirements of (1) continuity of the wave, and (2) no sudden changes in the slope of the wave. In other words, we assume that both the wave and its derivative are continuous functions.)



g / 1. A change in frequency without a change in wavelength would produce a discontinuity in the wave. 2. A simple change in wavelength without a reflection would result in a sharp kink in the wave.

Does this amount to a proof that reflection occurs? Not quite. We have only proved that certain types of wave motion are not valid solutions. In the following subsection, we prove that a valid solution can always be found in which a reflection occurs. Now in physics, we normally assume (but seldom prove formally) that the equations of motion have a unique solution, since otherwise a given set of initial conditions could lead to different behavior later on, but the Newtonian universe is supposed to be deterministic. Since the solution must be unique, and we derive below a valid solution involving a reflected pulse, we will have ended up with what amounts to a proof of reflection.

Intensity of reflection

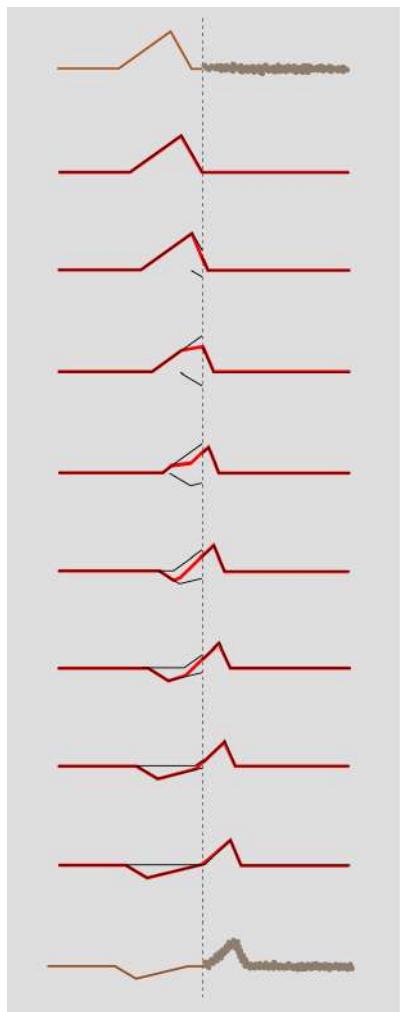
I will now show, in the case of waves on a string, that it is possible to satisfy the physical requirements given above by constructing a reflected wave, and as a bonus this will produce an equation for the proportions of reflection and transmission and a prediction as to which conditions will lead to inverted and which to uninverted reflection. We assume only that the principle of superposition holds, which is a good approximation for waves on a string of sufficiently small amplitude.

Let the unknown amplitudes of the reflected and transmitted waves be R and T , respectively. An inverted reflection would be represented by a negative value of R . We can without loss of generality take the incident (original) wave to have unit amplitude. Superposition tells us that if, for instance, the incident wave had double this amplitude, we could immediately find a corresponding solution simply by doubling R and T .

Just to the left of the boundary, the height of the wave is given by the height 1 of the incident wave, plus the height R of the part of the reflected wave that has just been created and begun heading back, for a total height of $1 + R$. On the right side immediately next to the boundary, the transmitted wave has a height T . To avoid a discontinuity, we must have

$$1 + R = T.$$

Next we turn to the requirement of equal slopes on both sides of the boundary. Let the slope of the incoming wave be s immediately to the left of the junction. If the wave was 100% reflected, and without inversion, then the slope of the reflected wave would be $-s$, since the wave has been reversed in direction. In general, the slope of the reflected wave equals $-sR$, and the slopes of the superposed waves on the left side add up to $s - sR$. On the right, the slope depends on the amplitude, T , but is also changed by the stretching or compression of the wave due to the change in speed. If, for example, the wave speed is twice as great on the right side, then the slope is cut in half by this effect. The slope on the right is



h / A pulse being partially reflected and partially transmitted at the boundary between two strings in which the wave speed is different. The top drawing shows the pulse heading to the right, toward the heavier string. For clarity, all but the first and last drawings are schematic. Once the reflected pulse begins to emerge from the boundary, it adds together with the trailing parts of the incident pulse. Their sum, shown as a wider line, is what is actually observed.

therefore $s(v_1/v_2)T$, where v_1 is the velocity in the original medium and v_2 the velocity in the new medium. Equality of slopes gives $s - sR = s(v_1/v_2)T$, or

$$1 - R = \frac{v_1}{v_2}T.$$

Solving the two equations for the unknowns R and T gives

$$R = \frac{v_2 - v_1}{v_2 + v_1}$$

and

$$T = \frac{2v_2}{v_2 + v_1}.$$

The first equation shows that there is no reflection unless the two wave speeds are different, and that the reflection is inverted in reflection back into a fast medium.

The energies of the transmitted and reflected waves always add up to the same as the energy of the original wave. There is never any abrupt loss (or gain) in energy when a wave crosses a boundary; conversion of wave energy to heat occurs for many types of waves, but it occurs throughout the medium. The equation for T , surprisingly, allows the amplitude of the transmitted wave to be greater than 1, i.e., greater than that of the incident wave. This does not violate conservation of energy, because this occurs when the second string is less massive, reducing its kinetic energy, and the transmitted pulse is broader and less strongly curved, which lessens its potential energy. In other words, the constant of proportionality in $E \propto A^2$ is different in the two different media.

We have attempted to develop some general facts about wave reflection by using the specific example of a wave on a string, which raises the question of whether these facts really are general. These issues are discussed in more detail in optional section 6.2.5, p. 389, but here is a brief summary.

The following facts are more generally true for wave reflection in one dimension.

- The wave is partially reflected and partially transmitted, with the reflected and transmitted parts sharing the energy.
- For an interface between media 1 and 2, there are two possible reflections: back into 1, and back into 2. One of these is inverting ($R < 0$) and the other is noninverting ($R > 0$).

The following aspects of our analysis may need to be modified for different types of waves.

- In some cases, the expressions for the reflected and transmitted amplitudes depend not on the ratio v_1/v_2 but on some more complicated ratio $v_1 \dots /v_2 \dots$, where \dots stands for some additional property of the medium.
- The sign of R , depends not just on this ratio but also on the type of the wave and on what we choose as a measure of amplitude.

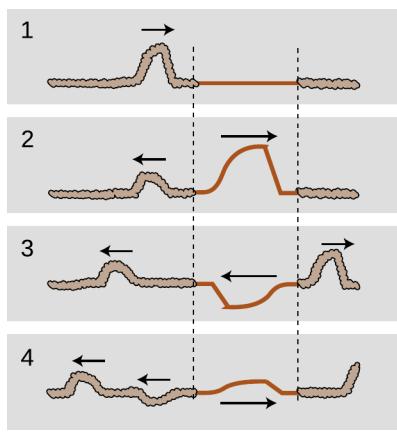
6.2.3 Interference effects

If you look at the front of a pair of high-quality binoculars, you will notice a greenish-blue coating on the lenses. This is advertised as a coating to prevent reflection. Now reflection is clearly undesirable — we want the light to go in the binoculars — but so far I've described reflection as an unalterable fact of nature, depending only on the properties of the two wave media. The coating can't change the speed of light in air or in glass, so how can it work? The key is that the coating itself is a wave medium. In other words, we have a three-layer sandwich of materials: air, coating, and glass. We will analyze the way the coating works, not because optical coatings are an important part of your education but because it provides a good example of the general phenomenon of wave interference effects.

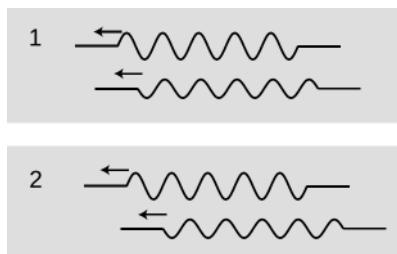
There are two different interfaces between media: an air-coating boundary and a coating-glass boundary. Partial reflection and partial transmission will occur at each boundary. For ease of visualization let's start by considering an equivalent system consisting of three dissimilar pieces of string tied together, and a wave pattern consisting initially of a single pulse. Figure i/1 shows the incident pulse moving through the heavy rope, in which its velocity is low. When it encounters the lighter-weight rope in the middle, a faster medium, it is partially reflected and partially transmitted. (The transmitted pulse is bigger, but nevertheless has only part of the original energy.) The pulse transmitted by the first interface is then partially reflected and partially transmitted by the second boundary, i/3. In figure i/4, two pulses are on the way back out to the left, and a single pulse is heading off to the right. (There is still a weak pulse caught between the two boundaries, and this will rattle back and forth, rapidly getting too weak to detect as it leaks energy to the outside with each partial reflection.)

Note how, of the two reflected pulses in i/4, one is inverted and one uninverted. One underwent reflection at the first boundary (a reflection back into a slower medium is uninverted), but the other was reflected at the second boundary (reflection back into a faster medium is inverted).

Now let's imagine what would have happened if the incoming wave pattern had been a long sinusoidal wave train instead of a single pulse. The first two waves to reemerge on the left could be



i / A pulse encounters two boundaries.



j / A sine wave has been reflected at two different boundaries, and the two reflections interfere.

in phase, $j/1$, or out of phase, $j/2$, or anywhere in between. The amount of lag between them depends entirely on the width of the middle segment of string. If we choose the width of the middle string segment correctly, then we can arrange for destructive interference to occur, $j/2$, with cancellation resulting in a very weak reflected wave.

This whole analysis applies directly to our original case of optical coatings. Visible light from most sources does consist of a stream of short sinusoidal wave-trains such as the ones drawn above. The only real difference between the waves-on-a-rope example and the case of an optical coating is that the first and third media are air and glass, in which light does not have the same speed. However, the general result is the same as long as the air and the glass have light-wave speeds that are either both greater than the coating's or both less than the coating's.

The business of optical coatings turns out to be a very arcane one, with a plethora of trade secrets and “black magic” techniques handed down from master to apprentice. Nevertheless, the ideas you have learned about waves in general are sufficient to allow you to come to some definite conclusions without any further technical knowledge. The self-check and discussion questions will direct you along these lines of thought.

self-check E

Color corresponds to wavelength of light waves. Is it possible to choose a thickness for an optical coating that will produce destructive interference for all colors of light?

▷ Answer, p. 1061

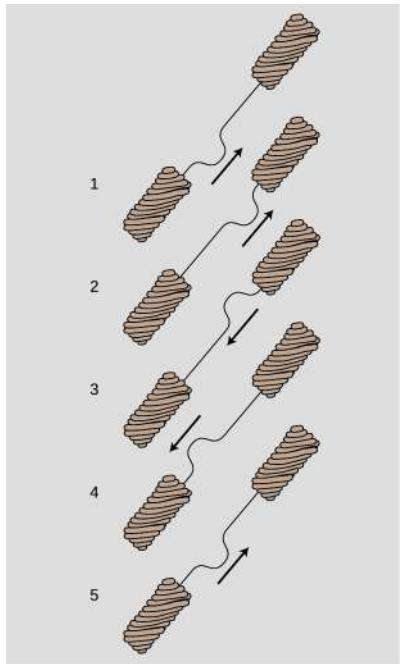
This example was typical of a wide variety of wave interference effects. With a little guidance, you are now ready to figure out for yourself other examples such as the rainbow pattern made by a compact disc or by a layer of oil on a puddle.

Discussion Questions

A Is it possible to get *complete* destructive interference in an optical coating, at least for light of one specific wavelength?

B Sunlight consists of sinusoidal wave-trains containing on the order of a hundred cycles back-to-back, for a length of something like a tenth of a millimeter. What happens if you try to make an optical coating thicker than this?

C Suppose you take two microscope slides and lay one on top of the other so that one of its edges is resting on the corresponding edge of the bottom one. If you insert a sliver of paper or a hair at the opposite end, a wedge-shaped layer of air will exist in the middle, with a thickness that changes gradually from one end to the other. What would you expect to see if the slides were illuminated from above by light of a single color? How would this change if you gradually lifted the lower edge of the top slide until the two slides were finally parallel?



k / A pulse bounces back and forth.



l / We model a guitar string attached to the guitar's body at both ends as a light-weight string attached to extremely heavy strings at its ends.

D An observation like the one described in discussion question C was used by Newton as evidence *against* the wave theory of light! If Newton didn't know about inverting and noninverting reflections, what would have seemed inexplicable to him about the region where the air layer had zero or nearly zero thickness?

6.2.4 Waves bounded on both sides

In the example of the previous section, it was theoretically true that a pulse would be trapped permanently in the middle medium, but that pulse was not central to our discussion, and in any case it was weakening severely with each partial reflection. Now consider a guitar string. At its ends it is tied to the body of the instrument itself, and since the body is very massive, the behavior of the waves when they reach the end of the string can be understood in the same way as if the actual guitar string was attached on the ends to strings that were extremely massive. Reflections are most intense when the two media are very dissimilar. Because the wave speed in the body is so radically different from the speed in the string, we should expect nearly 100% reflection.

Although this may seem like a rather bizarre physical model of the actual guitar string, it already tells us something interesting about the behavior of a guitar that we would not otherwise have understood. The body, far from being a passive frame for attaching the strings to, is actually the exit path for the wave energy in the strings. With every reflection, the wave pattern on the string loses a tiny fraction of its energy, which is then conducted through the body and out into the air. (The string has too little cross-section to make sound waves efficiently by itself.) By changing the properties of the body, moreover, we should expect to have an effect on the manner in which sound escapes from the instrument. This is clearly demonstrated by the electric guitar, which has an extremely massive, solid wooden body. Here the dissimilarity between the two wave media is even more pronounced, with the result that wave energy leaks out of the string even more slowly. This is why an electric guitar with no electric pickup can hardly be heard at all, and it is also the reason why notes on an electric guitar can be sustained for longer than notes on an acoustic guitar.

If we initially create a disturbance on a guitar string, how will the reflections behave? In reality, the finger or pick will give the string a triangular shape before letting it go, and we may think of this triangular shape as a very broad “dent” in the string which will spread out in both directions. For simplicity, however, let's just imagine a wave pattern that initially consists of a single, narrow pulse traveling up the neck, $k/1$. After reflection from the top end, it is inverted, $k/3$. Now something interesting happens: figure $k/5$ is identical to figure $k/1$. After two reflections, the pulse has been inverted twice and has changed direction twice. It is now back where it started. The motion is periodic. This is why a guitar produces

sounds that have a definite sensation of pitch.

self-check F

Notice that from $k/1$ to $k/5$, the pulse has passed by every point on the string exactly twice. This means that the total distance it has traveled equals $2L$, where L is the length of the string. Given this fact, what are the period and frequency of the sound it produces, expressed in terms of L and v , the velocity of the wave?

▷ Answer, p. 1061

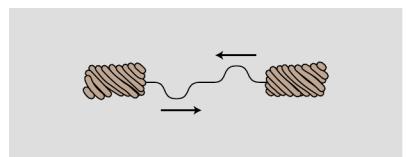
Note that if the waves on the string obey the principle of superposition, then the velocity must be independent of amplitude, and the guitar will produce the same pitch regardless of whether it is played loudly or softly. In reality, waves on a string obey the principle of superposition approximately, but not exactly. The guitar, like just about any acoustic instrument, is a little out of tune when played loudly. (The effect is more pronounced for wind instruments than for strings, but wind players are able to compensate for it.)

Now there is only one hole in our reasoning. Suppose we somehow arrange to have an initial setup consisting of two identical pulses heading toward each other, as in figure (g). They will pass through each other, undergo a single inverting reflection, and come back to a configuration in which their positions have been exactly interchanged. This means that the period of vibration is half as long. The frequency is twice as high.

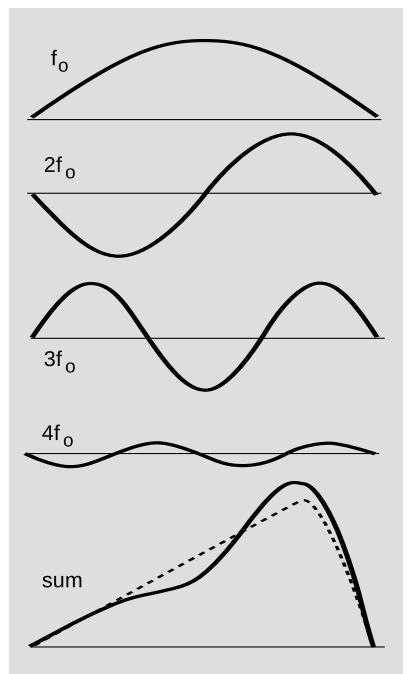
This might seem like a purely academic possibility, since nobody actually plays the guitar with two picks at once! But in fact it is an example of a very general fact about waves that are bounded on both sides. A mathematical theorem called Fourier's theorem states that any wave can be created by superposing sine waves. Figure n shows how even by using only four sine waves with appropriately chosen amplitudes, we can arrive at a sum which is a decent approximation to the realistic triangular shape of a guitar string being plucked. The one-hump wave, in which half a wavelength fits on the string, will behave like the single pulse we originally discussed. We call its frequency f_0 . The two-hump wave, with one whole wavelength, is very much like the two-pulse example. For the reasons discussed above, its frequency is $2f_0$. Similarly, the three-hump and four-hump waves have frequencies of $3f_0$ and $4f_0$.

Theoretically we would need to add together infinitely many such wave patterns to describe the initial triangular shape of the string exactly, although the amplitudes required for the very high frequency parts would be very small, and an excellent approximation could be achieved with as few as ten waves.

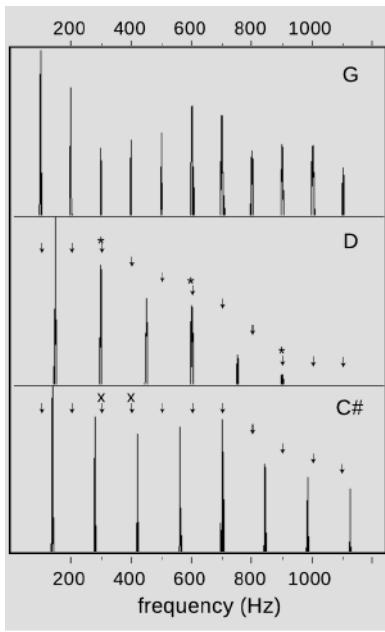
We thus arrive at the following very general conclusion. Whenever a wave pattern exists in a medium bounded on both sides by media in which the wave speed is very different, the motion can be broken down into the motion of a (theoretically infinite) series of sine



m / The period of this double-pulse pattern is half of what we'd otherwise expect.



n / Any wave can be made by superposing sine waves.



o / Graphs of loudness versus frequency for the vowel “ah,” sung as three different musical notes. G is consonant with D, since every overtone of G that is close to an overtone of D (marked “**”) is at exactly the same frequency. G and C \sharp are dissonant together, since some of the overtones of G (marked “x”) are close to, but not right on top of, those of C \sharp .

waves, with frequencies f_o , $2f_o$, $3f_o$, ... Except for some technical details, to be discussed below, this analysis applies to a vast range of sound-producing systems, including the air column within the human vocal tract. Because sounds composed of this kind of pattern of frequencies are so common, our ear-brain system has evolved so as to perceive them as a single, fused sensation of tone.

Musical applications

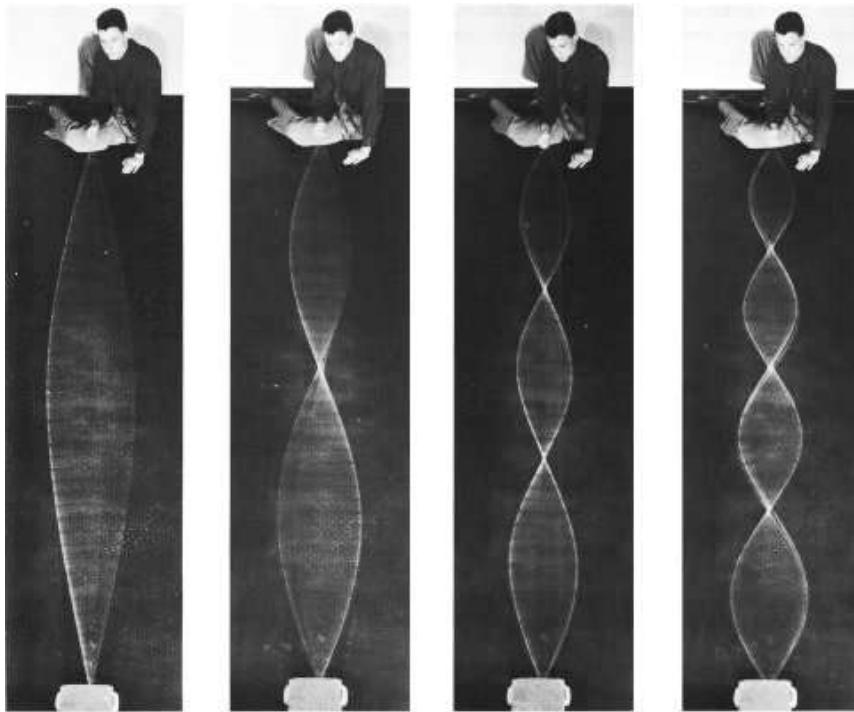
Many musicians claim to be able to pick out by ear several of the frequencies $2f_o$, $3f_o$, ..., called overtones or *harmonics* of the fundamental f_o , but they are kidding themselves. In reality, the overtone series has two important roles in music, neither of which depends on this fictitious ability to “hear out” the individual overtones.

First, the relative strengths of the overtones is an important part of the personality of a sound, called its timbre (rhymes with “amber”). The characteristic tone of the brass instruments, for example, is a sound that starts out with a very strong harmonic series extending up to very high frequencies, but whose higher harmonics die down drastically as the attack changes to the sustained portion of the note.

Second, although the ear cannot separate the individual harmonics of a single musical tone, it is very sensitive to clashes between the overtones of notes played simultaneously, i.e., in harmony. We tend to perceive a combination of notes as being dissonant if they have overtones that are close but not the same. Roughly speaking, strong overtones whose frequencies differ by more than 1% and less than 10% cause the notes to sound dissonant. It is important to realize that the term “dissonance” is not a negative one in music. No matter how long you search the radio dial, you will never hear more than three seconds of music without at least one dissonant combination of notes. Dissonance is a necessary ingredient in the creation of a musical cycle of tension and release. Musically knowledgeable people do not usually use the word “dissonant” as a criticism of music, and if they do, what they are really saying is that the dissonance has been used in a clumsy way, or without providing any contrast between dissonance and consonance.

Standing waves

Figure p shows sinusoidal wave patterns made by shaking a rope. I used to enjoy doing this at the bank with the pens on chains, back in the days when people actually went to the bank. You might think that I and the person in the photos had to practice for a long time in order to get such nice sine waves. In fact, a sine wave is the only shape that can create this kind of wave pattern, called a *standing wave*, which simply vibrates back and forth in one place without



p / Standing waves on a rope. (*PSSC Physics*.)

moving. The sine wave just creates itself automatically when you find the right frequency, because no other shape is possible.

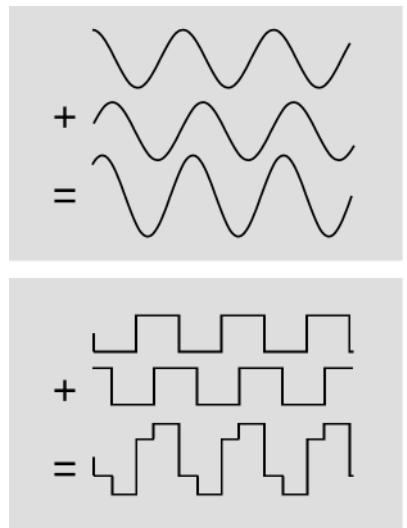
If you think about it, it's not even obvious that sine waves should be able to do this trick. After all, waves are supposed to travel at a set speed, aren't they? The speed isn't supposed to be zero! Well, we can actually think of a standing wave as a superposition of a moving sine wave with its own reflection, which is moving the opposite way. Sine waves have the unique mathematical property that the sum of sine waves of equal wavelength is simply a new sine wave with the same wavelength. As the two sine waves go back and forth, they always cancel perfectly at the ends, and their sum appears to stand still.

Standing wave patterns are rather important, since atoms are really standing-wave patterns of electron waves. You are a standing wave!

Standing-wave patterns of air columns

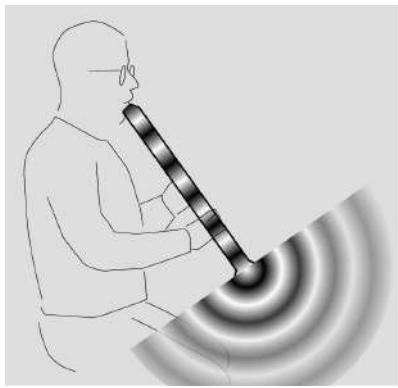
The air column inside a wind instrument behaves very much like the wave-on-a-string example we've been concentrating on so far, the main difference being that we may have either inverting or noninverting reflections at the ends.

Some organ pipes are closed at both ends. The speed of sound is different in metal than in air, so there is a strong reflection at the closed ends, and we can have standing waves. These reflections

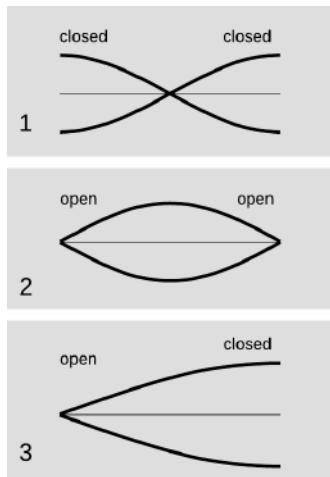


q / If you take a sine wave and make a copy of it shifted over, their sum is still a sine wave. The same is not true for a square wave.

are both density-noninverting, so we get symmetric standing-wave patterns, such as the one shown in figure s/1.



r / Surprisingly, sound waves undergo partial reflection at the open ends of tubes as well as closed ones.



s / Graphs of excess density versus position for the lowest-frequency standing waves of three types of air columns. Points on the axis have normal air density.

Figure r shows the sound waves in and around a bamboo Japanese flute called a shakuhachi, which is *open* at both ends of the air column. We can only have a standing wave pattern if there are reflections at the ends, but that is very counterintuitive — why is there any reflection at all, if the sound wave is free to emerge into open space, and there is no change in medium? Recall the reason why we got reflections at a change in medium: because the wavelength changes, so the wave has to readjust itself from one pattern to another, and the only way it can do that without developing a kink is if there is a reflection. Something similar is happening here. The only difference is that the wave is adjusting from being a plane wave to being a spherical wave. The reflections at the open ends are density-inverting, s/2, so the wave pattern is pinched off at the ends. Comparing panels 1 and 2 of the figure, we see that although the wave patterns are different, in both cases the wavelength is the same: in the lowest-frequency standing wave, half a wavelength fits inside the tube. Thus, it isn't necessary to memorize which type of reflection is inverting and which is uninverting. It's only necessary to know that the tubes are symmetric.

Finally, we can have an asymmetric tube: closed at one end and open at the other. A common example is the pan pipes, t, which are closed at the bottom and open at the top. The standing wave with the lowest frequency is therefore one in which 1/4 of a wavelength fits along the length of the tube, as shown in figure s/3.

Sometimes an instrument's physical appearance can be misleading. A concert flute, u, is closed at the mouth end and open at the other, so we would expect it to behave like an asymmetric air column; in reality, it behaves like a symmetric air column open at both ends, because the embouchure hole (the hole the player blows over) acts like an open end. The clarinet and the saxophone look similar, having a mouthpiece and reed at one end and an open end at the other, but they act different. In fact the clarinet's air column has patterns of vibration that are asymmetric, the saxophone symmetric. The discrepancy comes from the difference between the conical tube of the sax and the cylindrical tube of the clarinet. The adjustment of the wave pattern from a plane wave to a spherical wave is more gradual at the flaring bell of the saxophone.

self-check G

Draw a graph of pressure versus position for the first overtone of the air column in a tube open at one end and closed at the other. This will be the next-to-longest possible wavelength that allows for a point of maximum vibration at one end and a point of no vibration at the other. How many times shorter will its wavelength be compared to the frequency of the lowest-frequency standing wave, shown in the figure? Based on

this, how many times greater will its frequency be?
1061

▷ Answer, p.

The speed of sound

example 13

We can get a rough and ready derivation of the equation for the speed of sound by analyzing the standing waves in a cylindrical air column as a special type of Helmholtz resonance (example 25 on page 344), in which the cavity happens to have the same cross-sectional area as the neck. Roughly speaking, the regions of maximum density variation act like the cavity. The regions of minimum density variation, on the other hand, are the places where the velocity of the air is varying the most; these regions throttle back the speed of the vibration, because of the inertia of the moving air. If the cylinder has cross-sectional area A , then the “cavity” and “neck” parts of the wave both have lengths of something like $\lambda/2$, and the volume of the “cavity” is about $A\lambda/2$. We get $v = f\lambda = (\dots)\sqrt{\gamma P_0/\rho}$, where the factor (...) represents numerical stuff that we can’t possibly hope to have gotten right with such a crude argument. The correct result is in fact $v = \sqrt{\gamma P_0/\rho}$. Isaac Newton attempted the same calculation, but didn’t understand the thermodynamic effects involved, and therefore got a result that didn’t have the correct factor of $\sqrt{\gamma}$.

This chapter is summarized on page 1082. Notation and terminology are tabulated on pages 1070-1071.

6.2.5 * Some technical aspects of reflection

In this section we address some technical details of the treatment of reflection and transmission of waves.

Dependence of reflection on other variables besides velocity

In section 6.2.2 we derived the expressions for the transmitted and reflected amplitudes by demanding that two things match up on both sides of the boundary: the height of the wave and the slope of the wave. These requirements were stated purely in terms of kinematics (the description of how the wave moves) rather than dynamics (the explanation for the wave motion in terms of Newton’s laws). For this reason, the results depended only on the purely kinematic quantity $\alpha = v_2/v_1$, as can be seen more clearly if we rewrite the expressions in the following form:

$$R = \frac{\alpha - 1}{\alpha + 1} \quad \text{and} \quad T = \frac{2\alpha}{\alpha + 1}.$$

But this purely kinematical treatment only worked because of a dynamical fact that we didn’t emphasize. We assumed equality of the slopes, $s_1 = s_2$, because waves don’t like to have kinks. The underlying dynamical reason for this, in the case of a wave on a string, is that a kink is pointlike, so the portion of the string at the kink is infinitesimal in size, and therefore has essentially zero mass.



t / A pan pipe is an asymmetric air column, open at the top and closed at the bottom.



u / A concert flute looks like an asymmetric air column, open at the mouth end and closed at the other. However, its patterns of vibration are symmetric, because the embouchure hole acts like an open end.

If the transverse forces acting on it differed by some finite amount, then its acceleration would be infinite, which is not possible. The difference between the two forces is $Ts_1 - Ts_2$, so $s_1 = s_2$. But this relies on the assumption that T is the same on both sides of the boundary. Now this is true, because we can't put different amounts of tension on two ropes that are tied together end to end. Any excess tension applied to one tope is distributed equally to the other. For other types of waves, however, we cannot make a similar argument, and therefore it need not be true that $s_1 = s_2$.

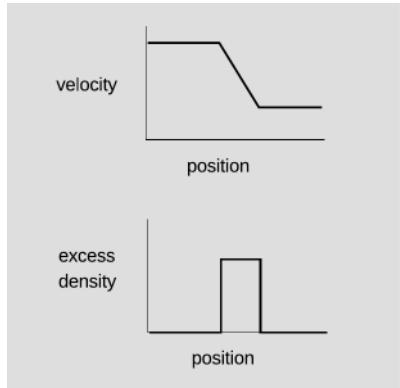
A more detailed analysis shows that in general we have not $\alpha = v_2/v_1$ but $\alpha = z_2/z_1$, where z is a quantity called impedance which is defined for this purpose. In a great many examples, as for the waves on a string, it is true that $v_2/v_1 = z_2/z_1$, but this is not a universal fact. Most of the exceptions are rather specialized and technical, such as the reflection of light waves when the media have magnetic properties, but one fairly common and important example is the case of sound waves, for which $z = \rho v$ depends not just on the wave velocity v but also on the density ρ . A practical example occurs in medical ultrasound scans, where the contrast of the image is made possible because of the very large differences in impedance between different types of tissue. The speed of sound in various tissues such as bone and muscle varies by about a factor of 2, which is not a particularly huge factor, but there are also large variations in density. The lung, for example, is basically a sponge or sack filled with air. For this reason, the acoustic impedances of the tissues show a huge amount of variation, with, e.g., $z_{\text{bone}}/z_{\text{lung}} \approx 40$.

Inverted and uninverted reflections in general

For waves on a string, reflections back into a faster medium are inverted, while those back into a slower medium are uninverted. Is this true for all types of waves? The rather subtle answer is that it depends on what property of the wave you are discussing.

Let's start by considering wave disturbances of freeway traffic. Anyone who has driven frequently on crowded freeways has observed the phenomenon in which one driver taps the brakes, starting a chain reaction that travels backward down the freeway as each person in turn exercises caution in order to avoid rear-ending anyone. The reason why this type of wave is relevant is that it gives a simple, easily visualized example of how our description of a wave depends on which aspect of the wave we have in mind. In steadily flowing freeway traffic, both the density of cars and their velocity are constant all along the road. Since there is no disturbance in this pattern of constant velocity and density, we say that there is no wave. Now if a wave is touched off by a person tapping the brakes, we can either describe it as a region of high density or as a region of decreasing velocity.

The freeway traffic wave is in fact a good model of a sound wave,



v / A disturbance in freeway traffic.



w / In the mirror image, the areas of positive excess traffic density are still positive, but the velocities of the cars have all been reversed, so areas of positive excess velocity have been turned into negative ones.

and a sound wave can likewise be described either by the density (or pressure) of the air or by its speed. Likewise many other types of waves can be described by either of two functions, one of which is often the derivative of the other with respect to position.

Now let's consider reflections. If we observe the freeway wave in a mirror, the high-density area will still appear high in density, but velocity in the opposite direction will now be described by a negative number. A person observing the mirror image will draw the same density graph, but the velocity graph will be flipped across the x axis, and its original region of negative slope will now have positive slope. Although I don't know any physical situation that would correspond to the reflection of a traffic wave, we can immediately apply the same reasoning to sound waves, which often do get reflected, and determine that a reflection can either be density-inverting and velocity-noninverting or density-noninverting and velocity-inverting.

This same type of situation will occur over and over as one encounters new types of waves, and to apply the analogy we need only determine which quantities, like velocity, become negated in a mirror image and which, like density, stay the same.

A light wave, for instance, consists of a traveling pattern of electric and magnetic fields. All you need to know in order to analyze the reflection of light waves is how electric and magnetic fields behave under reflection; you don't need to know any of the detailed physics of electricity and magnetism. An electric field can be detected, for example, by the way one's hair stands on end. The direction of the hair indicates the direction of the electric field. In a mirror image, the hair points the other way, so the electric field is apparently reversed in a mirror image. The behavior of magnetic fields, however, is a little tricky. The magnetic properties of a bar magnet, for instance, are caused by the aligned rotation of the outermost orbiting electrons of the atoms. In a mirror image, the direction of rotation is reversed, say from clockwise to counterclockwise, and so the magnetic field is reversed twice: once simply because the whole picture is flipped and once because of the reversed rotation of the electrons. In other words, magnetic fields do not reverse themselves in a mirror image. We can thus predict that there will be two possible types of reflection of light waves. In one, the electric field is inverted and the magnetic field uninverted (example 23, p. 728). In the other, the electric field is uninverted and the magnetic field inverted.

Problems

The symbols \checkmark , \blacksquare , etc. are explained on page 396.

- 1** The musical note middle C has a frequency of 262 Hz. What are its period and wavelength? \blacksquare

- 2** The following is a graph of the height of a water wave as a function of *position*, at a certain moment in time.



Trace this graph onto another piece of paper, and then sketch below it the corresponding graphs that would be obtained if

- (a) the amplitude and frequency were doubled while the velocity remained the same;
- (b) the frequency and velocity were both doubled while the amplitude remained unchanged;
- (c) the wavelength and amplitude were reduced by a factor of three while the velocity was doubled.

Explain all your answers. [Problem by Arnold Arons.] \blacksquare



Problem 3

- 3** (a) The graph shows the height of a water wave pulse as a function of position. Draw a graph of height as a function of time for a specific point on the water. Assume the pulse is traveling to the right.

- (b) Repeat part a, but assume the pulse is traveling to the left.
 - (c) Now assume the original graph was of height as a function of time, and draw a graph of height as a function of position, assuming the pulse is traveling to the right.
 - (d) Repeat part c, but assume the pulse is traveling to the left.
- Explain all your answers.* [Problem by Arnold Arons.] \blacksquare

- 4** At a particular moment in time, a wave on a string has a shape described by $y = 3.5 \cos(0.73\pi x + 0.45\pi t + 0.37\pi)$. The stuff inside the cosine is in radians. Assume that the units of the numerical constants are such that x , y , and t are in SI units. \triangleright Hint, p. 1036

- (a) Is the wave moving in the positive x or the negative x direction?
- (b) Find the wave's period, frequency, wavelength.
- (c) Find the wave's velocity.
- (d) Find the maximum velocity of any point on the string, and compare with the magnitude and direction of the wave's velocity.

\checkmark \blacksquare

5 The figure shows one wavelength of a steady sinusoidal wave traveling to the right along a string. Define a coordinate system in which the positive x axis points to the right and the positive y axis up, such that the flattened string would have $y = 0$. Copy the figure, and label with $y = 0$ all the appropriate parts of the string. Similarly, label with $v = 0$ all parts of the string whose velocities are zero, and with $a = 0$ all parts whose accelerations are zero. There is more than one point whose velocity is of the greatest magnitude. Pick one of these, and indicate the direction of its velocity vector. Do the same for a point having the maximum magnitude of acceleration. Explain all your answers.

[Problem by Arnold Arons.]



Problem 5.

6 (a) Find an equation for the relationship between the Doppler-shifted frequency of a wave and the frequency of the original wave, for the case of a stationary observer and a source moving directly toward or away from the observer. \checkmark

(b) Check that the units of your answer make sense.

(c) Check that the dependence on v_s makes sense.

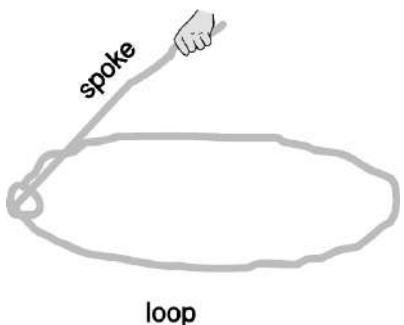
7 Suggest a quantitative experiment to look for any deviation from the principle of superposition for surface waves in water. Try to make your experiment simple and practical.

8 The simplest trick with a lasso is to spin a flat loop in a horizontal plane. The whirling loop of a lasso is kept under tension mainly due to its own rotation. Although the spoke's force on the loop has an inward component, we'll ignore it. The purpose of this problem, which is based on one by A.P. French, is to prove a cute fact about wave disturbances moving around the loop. As far as I know, this fact has no practical implications for trick roping! Let the loop have radius r and mass per unit length μ , and let its angular velocity be ω .

(a) Find the tension, T , in the loop in terms of r , μ , and ω . Assume the loop is a perfect circle, with no wave disturbances on it yet.

▷ Hint, p. 1036 ▷ Answer, p. 1069 \checkmark

(b) Find the velocity of a wave pulse traveling around the loop. Discuss what happens when the pulse moves in the same direction as the rotation, and when it travels contrary to the rotation. \checkmark



Problem 8.

9 A string hangs vertically, free at the bottom and attached at the top.

(a) Find the velocity of waves on the string as a function of the distance from the bottom. \checkmark

(b) Find the acceleration of waves on the string. ▷ Answer, p. 1069

(c) Interpret your answers to parts a and b for the case where a pulse comes down and reaches the end of the string. What happens next? Check your answer against experiment and conservation of energy.

10 Singing that is off-pitch by more than about 1% sounds bad. How fast would a singer have to be moving relative to the rest of a band to make this much of a change in pitch due to the Doppler effect? ✓ ■

11 Light travels faster in warmer air. On a sunny day, the sun can heat a road and create a layer of hot air above it. Let's model this layer as a uniform one with a sharp boundary separating it from the cooler air above. Use this model to explain the formation of a mirage appearing like the shiny surface of a pool of water. ■

12 (a) Compute the amplitude of light that is reflected back into air at an air-water interface, relative to the amplitude of the incident wave. Assume that the light arrives in the direction directly perpendicular to the surface. The speeds of light in air and water are 3.0×10^8 and 2.2×10^8 m/s, respectively.
(b) Find the energy of the reflected wave as a fraction of the incident energy.

▷ Hint, p. 1036 ✓ ■

13 A concert flute produces its lowest note, at about 262 Hz, when half of a wavelength fits inside its tube. Compute the length of the flute. ▷ Answer, p. 1069 ■

14 (a) A good tenor saxophone player can play all of the following notes without changing her fingering, simply by altering the tightness of her lips: E \flat (150 Hz), E \flat (300 Hz), B \flat (450 Hz), and E \flat (600 Hz). How is this possible? (I'm not asking you to analyze the coupling between the lips, the reed, the mouthpiece, and the air column, which is very complicated.)
(b) Some saxophone players are known for their ability to use this technique to play "freak notes," i.e., notes above the normal range of the instrument. Why isn't it possible to play notes below the normal range using this technique? ■

- C 261.6 Hz
- D 293.7
- E 329.6
- F 349.2
- G 392.0
- A 440.0
- B \flat 466.2

Problem 15.

15 The table gives the frequencies of the notes that make up the key of F major, starting from middle C and going up through all seven notes.

(a) Calculate the first four or five harmonics of C and G, and determine whether these two notes will be consonant or dissonant. (Recall that harmonics that differ by about 1-10% cause dissonance.)
(b) Do the same for C and B \flat . ■

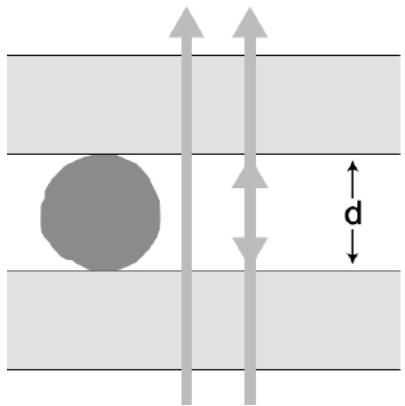
16 A Fabry-Perot interferometer, shown in the figure being used to measure the diameter of a thin filament, consists of two glass plates with an air gap between them. As the top plate is moved up or down with a screw, the light passing through the plates goes through a cycle of constructive and destructive interference, which is mainly due to interference between rays that pass straight through and those that are reflected twice back into the air gap. (Although the dimensions in this drawing are distorted for legibility, the glass plates would really be much thicker than the length of the wave-trains of light, so no interference effects would be observed due to reflections within the glass.)

- (a) If the top plate is cranked down so that the thickness, d , of the air gap is much less than the wavelength λ of the light, i.e., in the limit $d \rightarrow 0$, what is the phase relationship between the two rays? (Recall that the phase can be inverted by a reflection.) Is the interference constructive, or destructive?
- (b) If d is now slowly increased, what is the first value of d for which the interference is the same as at $d \rightarrow 0$? Express your answer in terms of λ .
- (c) Suppose the apparatus is first set up as shown in the figure. The filament is then removed, and n cycles of brightening and dimming are counted while the top plate is brought down to $d = 0$. What is the thickness of the filament, in terms of n and λ ?

Based on a problem by D.J. Raymond. ■

17 (a) A wave pulse moves into a new medium, where its velocity is greater by a factor α . Find an expression for the fraction, f , of the wave energy that is transmitted, in terms of α . Note that, as discussed in the text, you cannot simply find f by squaring the amplitude of the transmitted wave. ▷ Answer, p. 1069

(b) Suppose we wish to transmit a pulse from one medium to another, maximizing the fraction of the wave energy transmitted. To do so, we sandwich another layer in between them, so that the wave moves from the initial medium, where its velocity is v_1 , through the intermediate layer, where it is v_2 , and on into the final layer, where it becomes v_3 . What is the optimal value of v_2 ? (Assume that the middle layer is thicker than the length of the pulse, so there are no interference effects. Also, although there will be later echoes that are transmitted after multiple reflections back and forth across the middle layer, you are only to optimize the strength of the transmitted pulse that is first to emerge. In other words, it's simply a matter of applying your answer from part a twice to find the amount that finally gets through.) ▷ Answer, p. 1069 ■



Problem 16.

18 The expressions for the amplitudes of reflected and transmitted waves depend on the unitless ratio v_2/v_1 (or, more generally, on the ratio of the impedances). Call this ratio α . (a) Show that changing α to $1/\alpha$ (e.g., by interchanging the roles of the two media) has an effect on the reflected amplitude that can be expressed in a simple way, and discuss what this means in terms of inversion and energy. (b) Find the two values of α for which $|R| = 1/2$. ■

Key to symbols:

■ easy ■ typical ■ challenging ■ difficult ■ very difficult
✓ An answer check is available at www.lightandmatter.com.

Chapter 7

Relativity

7.1 Time is not absolute

When Einstein first began to develop the theory of relativity, around 1905, the only real-world observations he could draw on were ambiguous and indirect. Today, the evidence is part of everyday life. For example, every time you use a GPS receiver, a, you're using Einstein's theory of relativity. Somewhere between 1905 and today, technology became good enough to allow conceptually *simple* experiments that students in the early 20th century could only discuss in terms like "Imagine that we could..." A good jumping-on point is 1971. In that year, J.C. Hafele and R.E. Keating brought atomic clocks aboard commercial airliners, b, and went around the world, once from east to west and once from west to east. Hafele and Keating observed that there was a discrepancy between the times measured by the traveling clocks and the times measured by similar clocks that stayed home at the U.S. Naval Observatory in Washington. The east-going clock lost time, ending up off by -59 ± 10 nanoseconds, while the west-going one gained 273 ± 7 ns.

7.1.1 The correspondence principle

This establishes that time doesn't work the way Newton believed it did when he wrote that "Absolute, true, and mathematical time, of itself, and from its own nature flows equably without regard to anything external..." We are used to thinking of time as absolute and universal, so it is disturbing to find that it can flow at a different rate for observers in different frames of reference. Nevertheless, the effects that Hafele and Keating observed were small. This makes sense: Newton's laws have already been thoroughly tested by experiments under a wide variety of conditions, so a new theory like relativity must agree with Newton's to a good approximation, within the Newtonian theory's realm of applicability. This requirement of backward-compatibility is known as the correspondence principle.

7.1.2 Causality

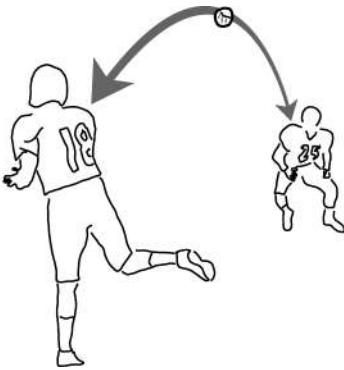
It's also reassuring that the effects on time were small compared to the three-day lengths of the plane trips. There was therefore no opportunity for paradoxical scenarios such as one in which the east-going experimenter arrived back in Washington before he left and then convinced himself not to take the trip. A theory that maintains



a / This Global Positioning System (GPS) system, running on a smartphone attached to a bike's handlebar, depends on Einstein's theory of relativity. Time flows at a different rate aboard a GPS satellite than it does on the bike, and the GPS software has to take this into account.



b / The clock took up two seats, and two tickets were bought for it under the name of "Mr. Clock."



c / Newton's laws do not distinguish past from future. The football could travel in either direction while obeying Newton's laws.

this kind of orderly relationship between cause and effect is said to satisfy causality.

Causality is like a water-hungry front-yard lawn in Los Angeles: we know we want it, but it's not easy to explain why. Even in plain old Newtonian physics, there is no clear distinction between past and future. In figure c, number 18 throws the football to number 25, and the ball obeys Newton's laws of motion. If we took a video of the pass and played it backward, we would see the ball flying from 25 to 18, and Newton's laws would still be satisfied. Nevertheless, we have a strong psychological impression that there is a forward arrow of time. I can remember what the stock market did last year, but I can't remember what it will do next year. Joan of Arc's military victories against England caused the English to burn her at the stake; it's hard to accept that Newton's laws provide an equally good description of a process in which her execution in 1431 caused her to win a battle in 1429. There is no consensus at this point among physicists on the origin and significance of time's arrow, and for our present purposes we don't need to solve this mystery. Instead, we merely note the empirical fact that, regardless of what causality really means and where it really comes from, its behavior is consistent. Specifically, experiments show that if an observer in a certain frame of reference observes that event A causes event B, then observers in other frames agree that A causes B, not the other way around. This is merely a generalization about a large body of experimental results, not a logically necessary assumption. If Keating had gone around the world and arrived back in Washington before he left, it would have disproved this statement about causality.

7.1.3 Time distortion arising from motion and gravity

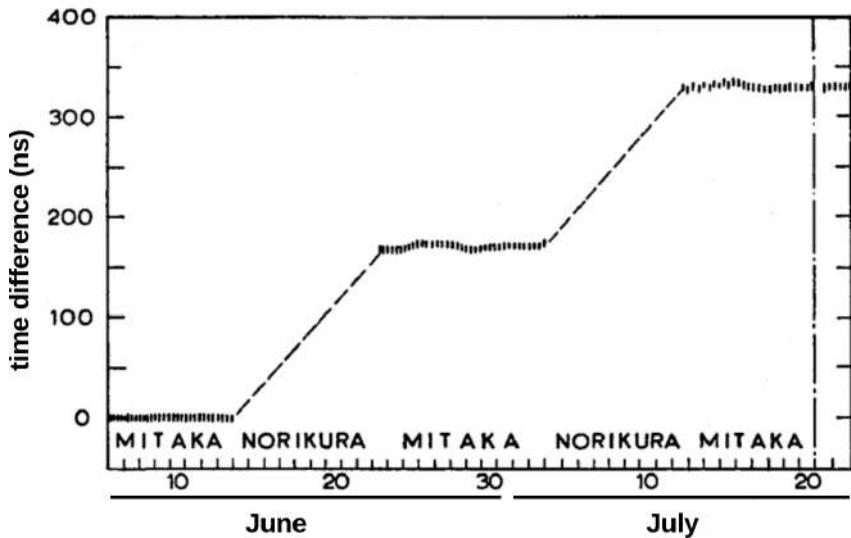
Hafele and Keating were testing specific quantitative predictions of relativity, and they verified them to within their experiment's error bars. Let's work backward instead, and inspect the empirical results for clues as to how time works.

The two traveling clocks experienced effects in opposite directions, and this suggests that the rate at which time flows depends on the motion of the observer. The east-going clock was moving in the same direction as the earth's rotation, so its velocity relative to the earth's center was greater than that of the clock that remained in Washington, while the west-going clock's velocity was correspondingly reduced. The fact that the east-going clock fell behind, and the west-going one got ahead, shows that the effect of motion is to make time go more slowly. This effect of motion on time was predicted by Einstein in his original 1905 paper on relativity, written when he was 26.

If this had been the only effect in the Hafele-Keating experiment, then we would have expected to see effects on the two flying clocks that were equal in size. Making up some simple numbers to keep the

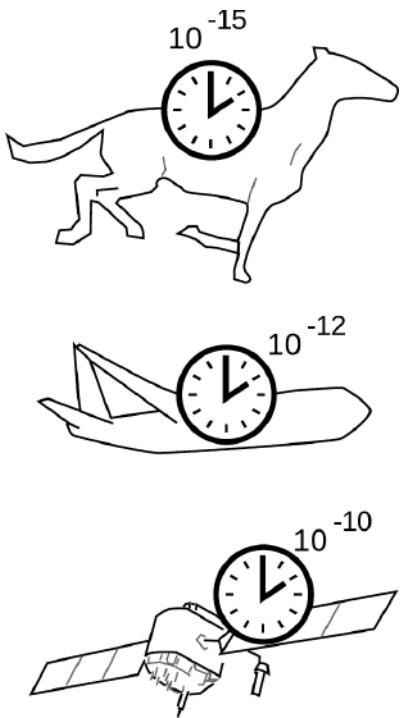
d / All three clocks are moving to the east. Even though the west-going plane is moving to the west relative to the air, the air is moving to the east due to the earth's rotation.

arithmetic transparent, suppose that the earth rotates from west to east at 1000 km/hr, and that the planes fly at 300 km/hr. Then the speed of the clock on the ground is 1000 km/hr, the speed of the clock on the east-going plane is 1300 km/hr, and that of the west-going clock 700 km/hr. Since the speeds of 700, 1000, and 1300 km/hr have equal spacing on either side of 1000, we would expect the discrepancies of the moving clocks relative to the one in the lab to be equal in size but opposite in sign.



e / A graph showing the time difference between two atomic clocks. One clock was kept at Mitaka Observatory, at 58 m above sea level. The other was moved back and forth to a second observatory, Norikura Corona Station, at the peak of the Norikura volcano, 2876 m above sea level. The plateaus on the graph are data from the periods when the clocks were compared side by side at Mitaka. The difference between one plateau and the next shows a gravitational effect on the rate of flow of time, accumulated during the period when the mobile clock was at the top of Norikura. Cf. problem 25, p. 462.

In fact, the two effects are unequal in size: -59 ns and 273 ns. This implies that there is a second effect involved, simply due to the planes' being up in the air. This was verified more directly in a 1978 experiment by Iijima and Fujiwara, figure e, in which identical atomic clocks were kept at rest at the top and bottom of a mountain near Tokyo. This experiment, unlike the Hafele-Keating one, isolates one effect on time, the gravitational one: time's rate of flow increases with height in a gravitational field. Einstein didn't figure out how to incorporate gravity into relativity until 1915, after much frustration and many false starts. The simpler version of the theory without gravity is known as special relativity, the full version as general relativity. We'll restrict ourselves to special relativity



f / The correspondence principle requires that the relativistic distortion of time become small for small velocities.

until section 7.4, and that means that what we want to focus on right now is the distortion of time due to motion, not gravity.

We can now see in more detail how to apply the correspondence principle. The behavior of the three clocks in the Hafele-Keating experiment shows that the amount of time distortion increases as the speed of the clock's motion increases. Newton lived in an era when the fastest mode of transportation was a galloping horse, and the best pendulum clocks would accumulate errors of perhaps a minute over the course of several days. A horse is much slower than a jet plane, so the distortion of time would have had a relative size of only $\sim 10^{-15}$ — much smaller than the clocks were capable of detecting. At the speed of a passenger jet, the effect is about 10^{-12} , and state-of-the-art atomic clocks in 1971 were capable of measuring that. A GPS satellite travels much faster than a jet airplane, and the effect on the satellite turns out to be $\sim 10^{-10}$. The general idea here is that all physical laws are approximations, and approximations aren't simply right or wrong in different situations. Approximations are better or worse in different situations, and the question is whether a particular approximation is good enough in a given situation to serve a particular purpose. The faster the motion, the worse the Newtonian approximation of absolute time. Whether the approximation is good enough depends on what you're trying to accomplish. The correspondence principle says that the approximation must have been good enough to explain all the experiments done in the centuries before Einstein came up with relativity.

By the way, don't get an inflated idea of the importance of the Hafele-Keating experiment. Special relativity had already been confirmed by a vast and varied body of experiments decades before 1971. The only reason I'm giving such a prominent role to this experiment, which was actually more important as a test of general relativity, is that it is conceptually very direct.

7.2 Distortion of space and time

7.2.1 The Lorentz transformation

Relativity says that when two observers are in different frames of reference, each observer considers the other one's perception of time to be distorted. We'll also see that something similar happens to their observations of distances, so both space and time are distorted. What exactly is this distortion? How do we even conceptualize it?

The idea isn't really as radical as it might seem at first. We can visualize the structure of space and time using a graph with position and time on its axes. These graphs are familiar by now, but we're going to look at them in a slightly different way. Before, we used them to describe the motion of objects. The grid underlying the graph was merely the stage on which the actors played their

parts. Now the background comes to the foreground: it's time and space themselves that we're studying. We don't necessarily need to have a line or a curve drawn on top of the grid to represent a particular object. We may, for example, just want to talk about events, depicted as points on the graph as in figure a. A distortion of the Cartesian grid underlying the graph can arise for perfectly ordinary reasons that Isaac Newton would have readily accepted. For example, we can simply change the units used to measure time and position, as in figure b.

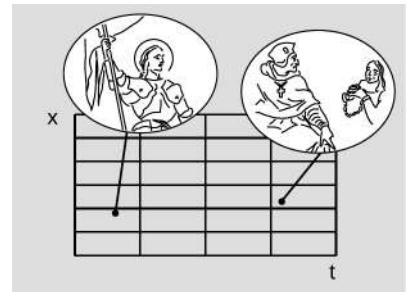
We're going to have quite a few examples of this type, so I'll adopt the convention shown in figure c for depicting them. Figure c summarizes the relationship between figures a and b in a more compact form. The gray rectangle represents the original coordinate grid of figure a, while the grid of black lines represents the new version from figure b. Omitting the grid from the gray rectangle makes the diagram easier to decode visually.

Our goal of unraveling the mysteries of special relativity amounts to nothing more than finding out how to draw a diagram like c in the case where the two different sets of coordinates represent measurements of time and space made by two different observers, each in motion relative to the other. Galileo and Newton thought they knew the answer to this question, but their answer turned out to be only approximately right. To avoid repeating the same mistakes, we need to clearly spell out what we think are the basic properties of time and space that will be a reliable foundation for our reasoning. I want to emphasize that there is no purely logical way of deciding on this list of properties. The ones I'll list are simply a summary of the patterns observed in the results from a large body of experiments. Furthermore, some of them are only approximate. For example, property 1 below is only a good approximation when the gravitational field is weak, so it is a property that applies to special relativity, not to general relativity.

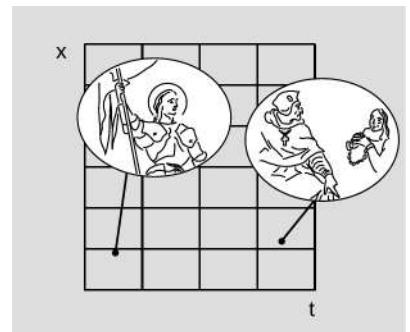
Experiments show that:

1. No point in time or space has properties that make it different from any other point.
2. Likewise, all directions in space have the same properties.
3. Motion is relative, i.e., all inertial frames of reference are equally valid.
4. Causality holds, in the sense described on page 397.
5. Time depends on the state of motion of the observer.

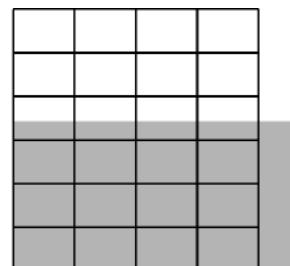
Most of these are not very subversive. Properties 1 and 2 date back to the time when Galileo and Newton started applying the



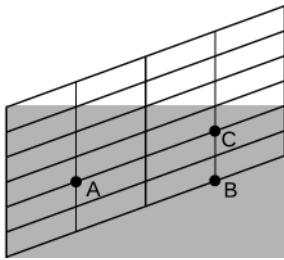
a / Two events are given as points on a graph of position versus time. Joan of Arc helps to restore Charles VII to the throne. At a later time and a different position, Joan of Arc is sentenced to death.



b / A change of units distorts an x - t graph. This graph depicts exactly the same events as figure a. The only change is that the x and t coordinates are measured using different units, so the grid is compressed in t and expanded in x .



c / A convention we'll use to represent a distortion of time and space.



d / A Galilean version of the relationship between two frames of reference. As in all such graphs in this chapter, the original coordinates, represented by the gray rectangle, have a time axis that goes to the right, and a position axis that goes straight up.

same universal laws of motion to the solar system and to the earth; this contradicted Aristotle, who believed that, for example, a rock would naturally want to move in a certain special direction (down) in order to reach a certain special location (the earth's surface). Property 3 is the reason that Einstein called his theory "relativity," but Galileo and Newton believed exactly the same thing to be true, as dramatized by Galileo's run-in with the Church over the question of whether the earth could really be in motion around the sun. Property 4 would probably surprise most people only because it asserts in such a weak and specialized way something that they feel deeply must be true. The only really strange item on the list is 5, but the Hafele-Keating experiment forces it upon us.

If it were not for property 5, we could imagine that figure d would give the correct transformation between frames of reference in motion relative to one another. Let's say that observer 1, whose grid coincides with the gray rectangle, is a hitch-hiker standing by the side of a road. Event A is a raindrop hitting his head, and event B is another raindrop hitting his head. He says that A and B occur at the same location in space. Observer 2 is a motorist who drives by without stopping; to him, the passenger compartment of his car is at rest, while the asphalt slides by underneath. He says that A and B occur at different points in space, because during the time between the first raindrop and the second, the hitch-hiker has moved backward. On the other hand, observer 2 says that events A and C occur in the same place, while the hitch-hiker disagrees. The slope of the grid-lines is simply the velocity of the relative motion of each observer relative to the other.

Figure d has familiar, comforting, and eminently sensible behavior, but it also happens to be wrong, because it violates property 5. The distortion of the coordinate grid has only moved the vertical lines up and down, so both observers agree that events like B and C are simultaneous. If this was really the way things worked, then all observers could synchronize all their clocks with one another for once and for all, and the clocks would never get out of sync. This contradicts the results of the Hafele-Keating experiment, in which all three clocks were initially synchronized in Washington, but later went out of sync because of their different states of motion.

It might seem as though we still had a huge amount of wiggle room available for the correct form of the distortion. It turns out, however, that properties 1-5 are sufficient to prove that there is only one answer, which is the one found by Einstein in 1905. To see why this is, let's work by a process of elimination.

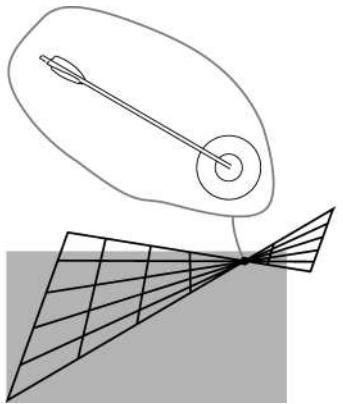
Figure e shows a transformation that might seem at first glance to be as good a candidate as any other, but it violates property 3, that motion is relative, for the following reason. In observer 2's frame of reference, some of the grid lines cross one another. This

means that observers 1 and 2 disagree on whether or not certain events are the same. For instance, suppose that event A marks the arrival of an arrow at the bull's-eye of a target, and event B is the location and time when the bull's-eye is punctured. Events A and B occur at the same location and at the same time. If one observer says that A and B coincide, but another says that they don't, we have a direct contradiction. Since the two frames of reference in figure e give contradictory results, one of them is right and one is wrong. This violates property 3, because all inertial frames of reference are supposed to be equally valid. To avoid problems like this, we clearly need to make sure that none of the grid lines ever cross one another.

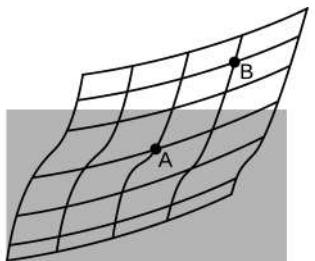
The next type of transformation we want to kill off is shown in figure f, in which the grid lines curve, but never cross one another. The trouble with this one is that it violates property 1, the uniformity of time and space. The transformation is unusually “twisty” at A, whereas at B it’s much more smooth. This can’t be correct, because the transformation is only supposed to depend on the relative state of motion of the two frames of reference, and that given information doesn’t single out a special role for any particular point in spacetime. If, for example, we had one frame of reference *rotating* relative to the other, then there would be something special about the axis of rotation. But we’re only talking about *inertial* frames of reference here, as specified in property 3, so we can’t have rotation; each frame of reference has to be moving in a straight line at constant speed. For frames related in this way, there is nothing that could single out an event like A for special treatment compared to B, so transformation f violates property 1.

The examples in figures e and f show that the transformation we’re looking for must be linear, meaning that it must transform lines into lines, and furthermore that it has to take parallel lines to parallel lines. Einstein wrote in his 1905 paper that “...on account of the property of homogeneity [property 1] which we ascribe to time and space, the [transformation] must be linear.”¹ Applying this to our diagrams, the original gray rectangle, which is a special type of parallelogram containing right angles, must be transformed into another parallelogram. There are three types of transformations, figure g, that have this property. Case I is the Galilean transformation of figure d on page 402, which we’ve already ruled out.

Case II can also be discarded. Here every point on the grid rotates counterclockwise. What physical parameter would determine the amount of rotation? The only thing that could be relevant would be v , the relative velocity of the motion of the two frames of reference with respect to one another. But if the angle of rotation was pro-

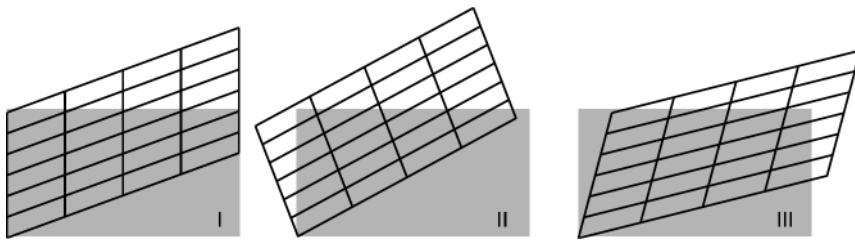


e / A transformation that leads to disagreements about whether two events occur at the same time and place. This is not just a matter of opinion. Either the arrow hit the bull's-eye or it didn't.

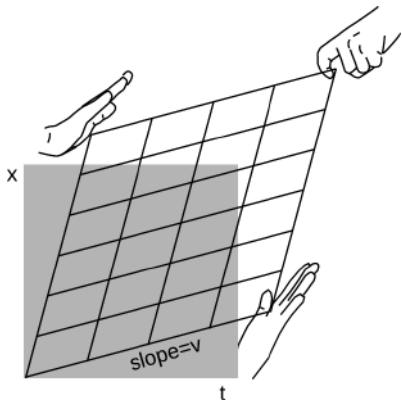


f / A nonlinear transformation.

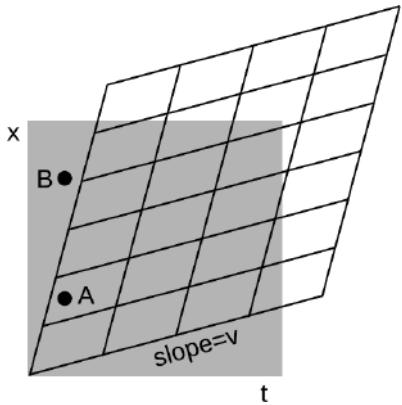
¹A. Einstein, “On the Electrodynamics of Moving Bodies,” *Annalen der Physik* 17 (1905), p. 891, tr. Saha and Bose.



g / Three types of transformations that preserve parallelism. Their distinguishing feature is what they do to simultaneity, as shown by what happens to the left edge of the original rectangle. In I, the left edge remains vertical, so simultaneous events remain simultaneous. In II, the left edge turns counterclockwise. In III, it turns clockwise.



h / In the units that are most convenient for relativity, the transformation has symmetry about a 45-degree diagonal line.



i / Interpretation of the Lorentz transformation. The slope indicated in the figure gives the relative velocity of the two frames of reference. Events A and B that were simultaneous in frame 1 are not simultaneous in frame 2, where event A occurs to the right of the $t = 0$ line represented by the left edge of the grid, but event B occurs to its left.

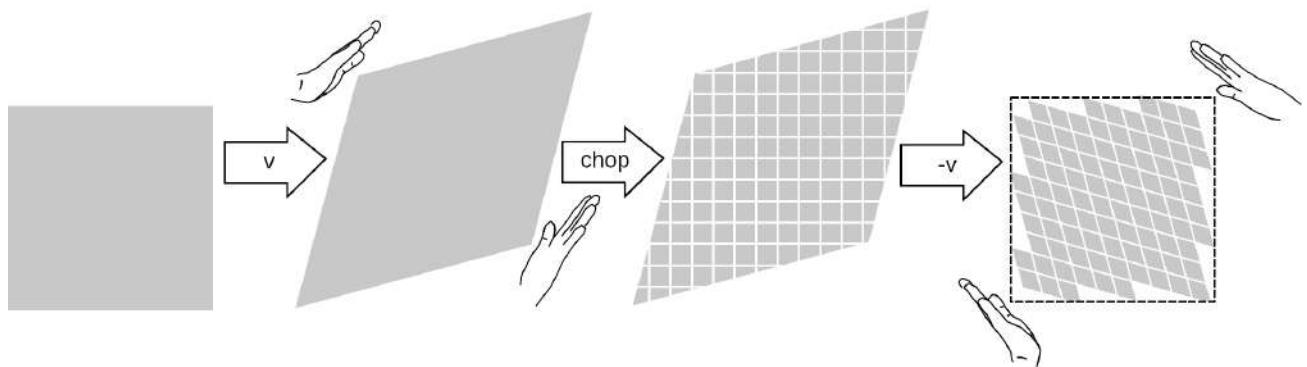
portional to v , then for large enough velocities the grid would have left and right reversed, and this would violate property 4, causality: one observer would say that event A caused a later event B, but another observer would say that B came first and caused A.

The only remaining possibility is case III, which I've redrawn in figure h with a couple of changes. This is the one that Einstein predicted in 1905. The transformation is known as the Lorentz transformation, after Hendrik Lorentz (1853-1928), who partially anticipated Einstein's work, without arriving at the correct interpretation. The distortion is a kind of smooshing and stretching, as suggested by the hands. Also, we've already seen in figures a-c on page 401 that we're free to stretch or compress everything as much as we like in the horizontal and vertical directions, because this simply corresponds to choosing different units of measurement for time and distance. In figure h I've chosen units that give the whole drawing a convenient symmetry about a 45-degree diagonal line. Ordinarily it wouldn't make sense to talk about a 45-degree angle on a graph whose axes had different units. But in relativity, the symmetric appearance of the transformation tells us that space and time ought to be treated on the same footing, and measured in the same units.

As in our discussion of the Galilean transformation, slopes are interpreted as velocities, and the slope of the near-horizontal lines in figure i is interpreted as the relative velocity of the two observers. The difference between the Galilean version and the relativistic one is that now there is smooshing happening from the other side as well. Lines that were vertical in the original grid, representing simultaneous events, now slant over to the right. This tells us that, as required by property 5, different observers do not agree on whether events that occur in different places are simultaneous. The Hafele-Keating experiment tells us that this non-simultaneity effect is fairly small, even when the velocity is as big as that of a passenger jet, and this is what we would have anticipated by the correspondence principle. The way that this is expressed in the graph is that if we

pick the time unit to be the second, then the distance unit turns out to be hundreds of thousands of miles. In these units, the velocity of a passenger jet is an extremely small number, so the slope v in figure i is extremely small, and the amount of distortion is tiny — it would be much too small to see on this scale.

The only thing left to determine about the Lorentz transformation is the size of the transformed parallelogram relative to the size of the original one. Although the drawing of the hands in figure h may suggest that the grid deforms like a framework made of rigid coat-hanger wire, that is not the case. If you look carefully at the figure, you'll see that the edges of the smooshed parallelogram are actually a little longer than the edges of the original rectangle. In fact what stays the same is not lengths but *areas*, as proved in the caption to figure j.



j / Proof that Lorentz transformations don't change area: We first subject a square to a transformation with velocity v , and this increases its area by a factor $R(v)$, which we want to prove equals 1. We chop the resulting parallelogram up into little squares and finally apply a $-v$ transformation; this changes each little square's area by a factor $R(-v)$, so the whole figure's area is also scaled by $R(-v)$. The final result is to restore the square to its original shape and area, so $R(v)R(-v) = 1$. But $R(v) = R(-v)$ by property 2 of spacetime on page 401, which states that all directions in space have the same properties, so $R(v) = 1$.

7.2.2 The γ factor

With a little algebra and geometry (homework problem 7, page 458), one can use the equal-area property to show that the factor γ (Greek letter gamma) defined in figure k is given by the equation

$$\gamma = \frac{1}{\sqrt{1 - v^2}}.$$

If you've had good training in physics, the first thing you probably think when you look at this equation is that it must be nonsense, because its units don't make sense. How can we take something with units of velocity squared, and subtract it from a unitless 1? But remember that this is expressed in our special relativistic units, in which the same units are used for distance and time. We refer

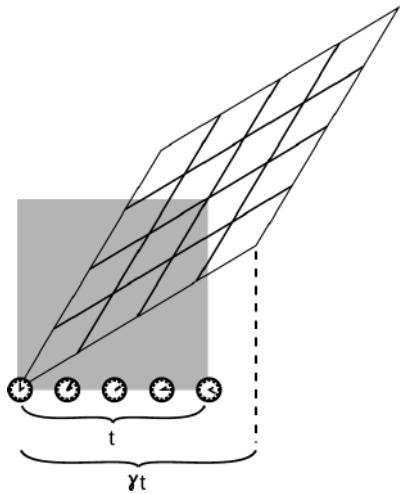
to these as *natural* units. In this system, velocities are always unitless. This sort of thing happens frequently in physics. For instance, before James Joule discovered conservation of energy, nobody knew that heat and mechanical energy were different forms of the same thing, so instead of measuring them both in units of joules as we would do now, they measured heat in one unit (such as calories) and mechanical energy in another (such as foot-pounds). In ordinary metric units, we just need an extra conversion factor c , and the equation becomes

$$\gamma = \frac{1}{\sqrt{1 - (\frac{v}{c})^2}}.$$

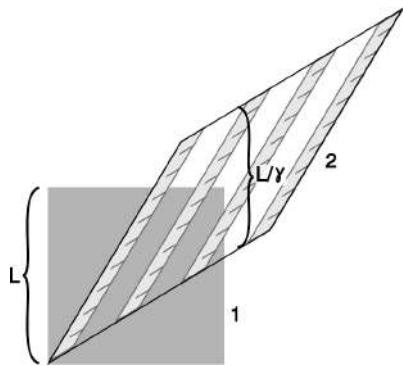
Here's why we care about γ . Figure k defines it as the ratio of two times: the time between two events as expressed in one coordinate system, and the time between the same two events as measured in the other one. The interpretation is:

Time dilation

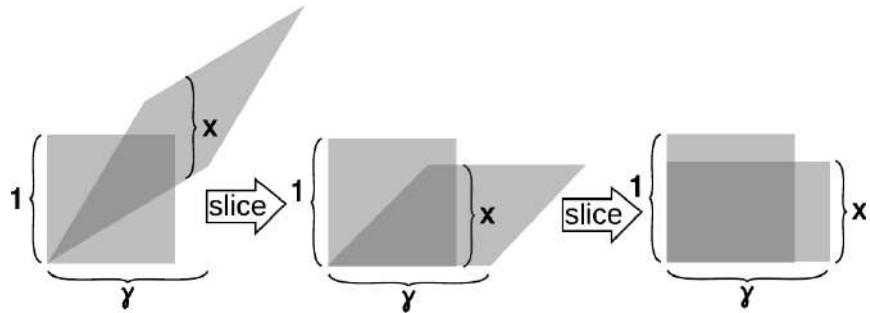
A clock runs fastest in the frame of reference of an observer who is at rest relative to the clock. An observer in motion relative to the clock at speed v perceives the clock as running more slowly by a factor of γ .



k / The γ factor.



l / The ruler is moving in frame 1, represented by a square, but at rest in frame 2, shown as a parallelogram. Each picture of the ruler is a snapshot taken at a certain moment as judged according to frame 2's notion of simultaneity. An observer in frame 1 judges the ruler's length instead according to frame 1's definition of simultaneity, i.e., using points that are lined up vertically on the graph. The ruler appears shorter in the frame in which it is moving. As proved in figure m, the length contracts from L to L/γ .



m / This figure proves, as claimed in figure l, that the length contraction is $x = 1/\gamma$. First we slice the parallelogram vertically like a salami and slide the slices down, making the top and bottom edges horizontal. Then we do the same in the horizontal direction, forming a rectangle with sides γ and x . Since both the Lorentz transformation and the slicing processes leave areas unchanged, the area γx of the rectangle must equal the area of the original square, which is 1.

As proved in figures l and m, lengths are also distorted:

Length contraction

A meter-stick appears longest to an observer who is at rest relative to it. An observer moving relative to the meter-stick at v observes the stick to be shortened by a factor of γ .

self-check A

What is γ when $v = 0$? What does this mean?

▷ Answer, p. 1062

Figure n shows the behavior of γ as a function of v .

Changing an equation from natural units to SI *example 1*

Often it is easier to do all of our algebra in natural units, which are simpler because $c = 1$, and all factors of c can therefore be omitted. For example, suppose we want to solve for v in terms of γ . In natural units, we have $\gamma = 1/\sqrt{1 - v^2}$, so $\gamma^{-2} = 1 - v^2$, and $v = \sqrt{1 - \gamma^{-2}}$.

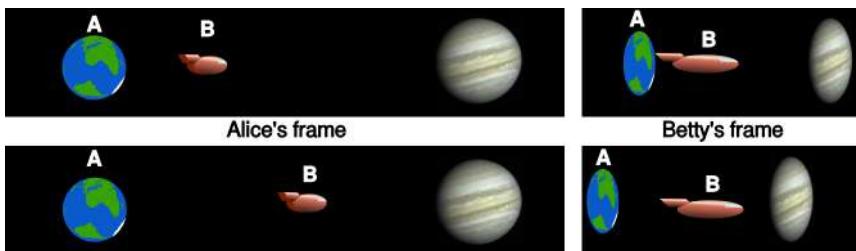
This form of the result might be fine for many purposes, but if we wanted to find a value of v in SI units, we would need to reinsert factors of c in the final result. There is no need to do this throughout the whole derivation. By looking at the final result, we see that there is only one possible way to do this so that the results make sense in SI, which is to write $v = c\sqrt{1 - \gamma^{-2}}$.

Motion of a ray of light *example 2*

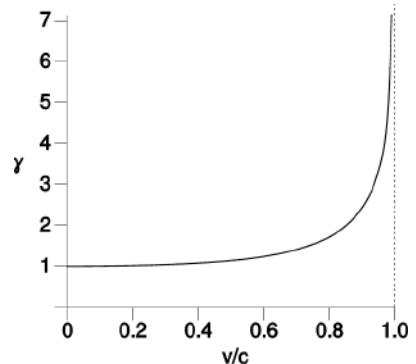
- ▷ The motion of a certain ray of light is given by the equation $x = -t$. Is this expressed in natural units, or in SI units? Convert to the other system.
- ▷ The equation is in natural units. It wouldn't make sense in SI units, because we would have meters on the left and seconds on the right. To convert to SI units, we insert a factor of c in the only possible place that will cause the equation to make sense: $x = -ct$.

An interstellar road trip *example 3*

Alice stays on earth while her twin Betty heads off in a spaceship for Tau Ceti, a nearby star. Tau Ceti is 12 light-years away, so even though Betty travels at 87% of the speed of light, it will take her a long time to get there: 14 years, according to Alice.



Betty experiences time dilation. At this speed, her γ is 2.0, so that the voyage will only seem to her to last 7 years. But there is perfect symmetry between Alice's and Betty's frames of reference, so Betty agrees with Alice on their relative speed; Betty sees herself as being at rest, while the sun and Tau Ceti both move backward at 87% of the speed of light. How, then, can she observe Tau Ceti to get to her in only 7 years, when it should take 14 years to travel 12 light-years at this speed?



n / A graph of γ as a function of v .

o / Example 3.

We need to take into account length contraction. Betty sees the distance between the sun and Tau Ceti to be shrunk by a factor of 2. The same thing occurs for Alice, who observes Betty and her spaceship to be foreshortened.

The correspondence principle

example 4

The correspondence principle requires that γ be close to 1 for the velocities much less than c encountered in everyday life. In natural units, $\gamma = (1 - v^2)^{-1/2}$. For small values of v , the approximation $(1 + \epsilon)^{\rho} \approx 1 + \rho\epsilon$ holds (see p. 1022). Applying this approximation, we find $\gamma \approx 1 + v^2/2$.

As expected, this gives approximately 1 when v is small compared to 1 (i.e., compared to c , which equals 1 in natural units).

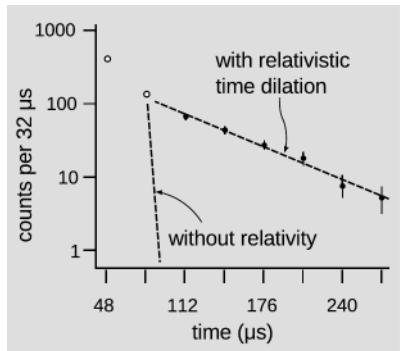
In problem 32 on p. 464 we rewrite this in SI units.

Figure n on p. 407 shows that the approximation is *not* valid for large values of v/c . In fact, γ blows up to infinity as v gets closer and closer to c .

Large time dilation

example 5

The time dilation effect in the Hafele-Keating experiment was very small. If we want to see a large time dilation effect, we can't do it with something the size of the atomic clocks they used; the kinetic energy would be greater than the total megatonnage of all the world's nuclear arsenals. We can, however, accelerate subatomic particles to speeds at which γ is large. For experimental particle physicists, relativity is something you do all day before heading home and stopping off at the store for milk. An early, low-precision experiment of this kind was performed by Rossi and Hall in 1941, using naturally occurring cosmic rays. Figure q shows a 1974 experiment² of a similar type which verified the time dilation predicted by relativity to a precision of about one part per thousand.



p / Muons accelerated to nearly c undergo radioactive decay much more slowly than they would according to an observer at rest with respect to the muons. The first two data-points (unfilled circles) were subject to large systematic errors.

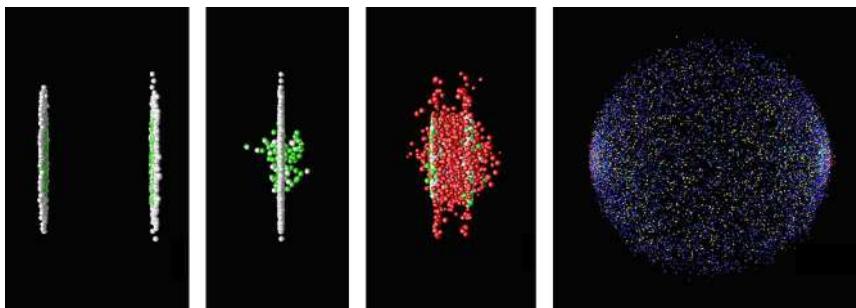
Particles called muons (named after the Greek letter μ , “myoo”) were produced by an accelerator at CERN, near Geneva. A muon is essentially a heavier version of the electron. Muons undergo radioactive decay, lasting an average of only $2.197 \mu s$ before they evaporate into an electron and two neutrinos. The 1974 experiment was actually built in order to measure the magnetic properties of muons, but it produced a high-precision test of time dilation as a byproduct. Because muons have the same electric charge as electrons, they can be trapped using magnetic fields. Muons were injected into the ring shown in figure q, circling around it until they underwent radioactive decay. At the speed at which these muons were traveling, they had $\gamma = 29.33$, so on the average they

²Bailey et al., Nucl. Phys. B150(1979) 1



q / Apparatus used for the test of relativistic time dilation described in example 5. The prominent black and white blocks are large magnets surrounding a circular pipe with a vacuum inside. (c) 1974 by CERN.

lasted 29.33 times longer than the normal lifetime. In other words, they were like tiny alarm clocks that self-destructed at a randomly selected time. Figure p shows the number of radioactive decays counted, as a function of the time elapsed after a given stream of muons was injected into the storage ring. The two dashed lines show the rates of decay predicted with and without relativity. The relativistic line is the one that agrees with experiment.

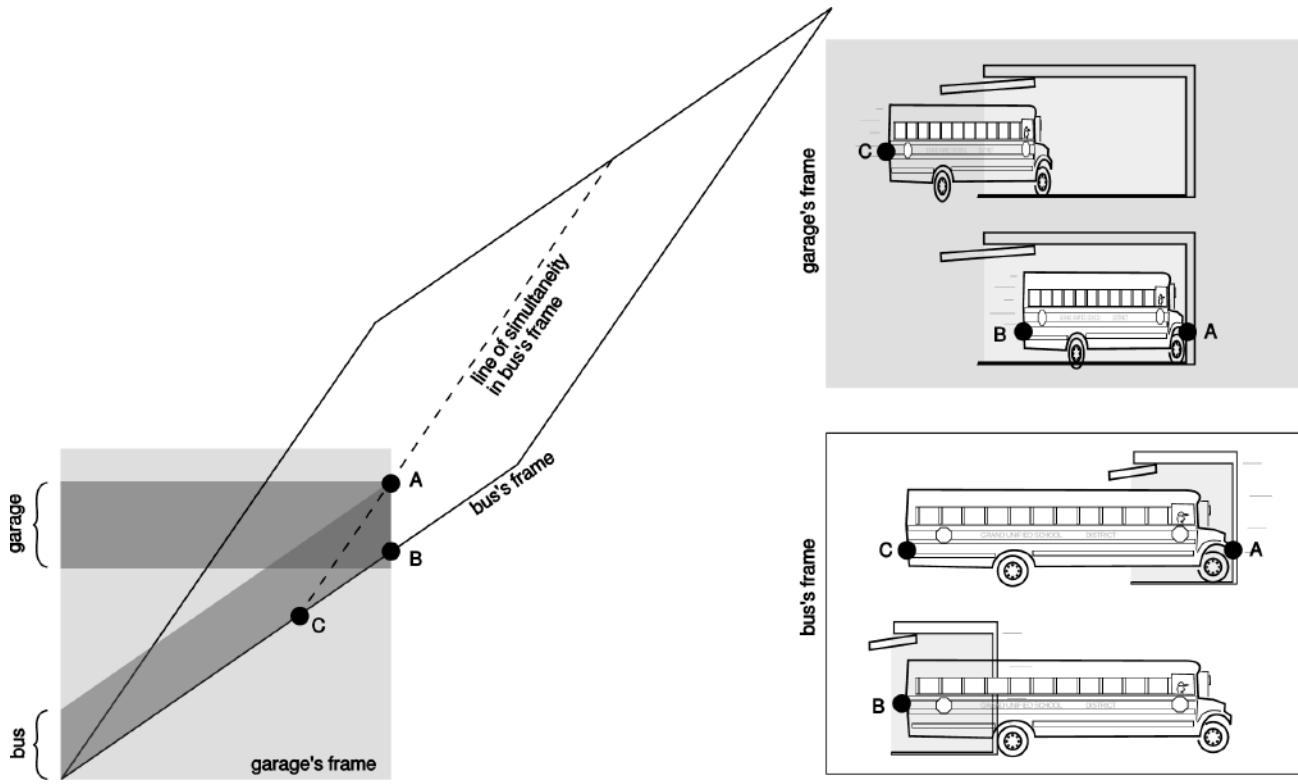


r / Colliding nuclei show relativistic length contraction.

An example of length contraction

Figure r shows an artist's rendering of the length contraction for the collision of two gold nuclei at relativistic speeds in the RHIC accelerator in Long Island, New York. The gold nuclei would appear nearly spherical (or just slightly lengthened like an American football) in frames moving along with them, but in the laboratory's frame, they both appear drastically foreshortened as they approach the point of collision. The later pictures show the nuclei merging to form a hot soup, observed at RHIC in 2010, in which the quarks are no longer confined inside the protons and neutrons.

example 6



s / Example 7: In the garage's frame of reference, the bus is moving, and can fit in the garage due to its length contraction. In the bus's frame of reference, the garage is moving, and can't hold the bus due to *its* length contraction.

The garage paradox

example 7

One of the most famous of all the so-called relativity paradoxes has to do with our incorrect feeling that simultaneity is well defined. The idea is that one could take a schoolbus and drive it at relativistic speeds into a garage of ordinary size, in which it normally would not fit. Because of the length contraction, the bus would supposedly fit in the garage. The driver, however, will perceive the *garage* as being contracted and thus even less able to contain the bus.

The paradox is resolved when we recognize that the concept of fitting the bus in the garage “all at once” contains a hidden assumption, the assumption that it makes sense to ask whether the front and back of the bus can *simultaneously* be in the garage. Observers in different frames of reference moving at high relative speeds do not necessarily agree on whether things happen simultaneously. As shown in figure s, the person in the garage's frame can shut the door at an instant B he perceives to be simultaneous with the front bumper's arrival A at the back wall of the garage, but the driver would not agree about the simultaneity of these two events, and would perceive the door as having shut long after she plowed through the back wall.

7.2.3 The universal speed c

Let's think a little more about the role of the 45-degree diagonal in the Lorentz transformation. Slopes on these graphs are interpreted as velocities. This line has a slope of 1 in relativistic units, but that slope corresponds to c in ordinary metric units. We already know that the relativistic distance unit must be extremely large compared to the relativistic time unit, so c must be extremely large. Now note what happens when we perform a Lorentz transformation: this particular line gets stretched, but the new version of the line lies right on top of the old one, and its slope stays the same. In other words, if one observer says that something has a velocity equal to c , every other observer will agree on that velocity as well. (The same thing happens with $-c$.)

Velocities don't simply add and subtract.

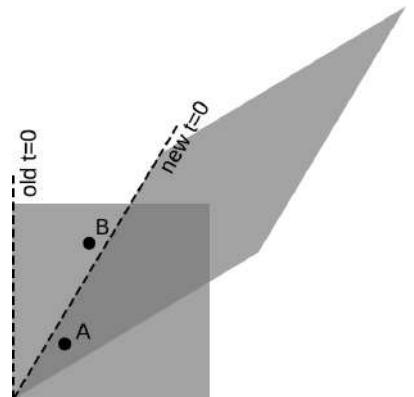
This is counterintuitive, since we expect velocities to add and subtract in relative motion. If a dog is running away from me at 5 m/s relative to the sidewalk, and I run after it at 3 m/s, the dog's velocity in my frame of reference is 2 m/s. According to everything we have learned about motion, the dog must have different speeds in the two frames: 5 m/s in the sidewalk's frame and 2 m/s in mine. But velocities are measured by dividing a distance by a time, and both distance and time are distorted by relativistic effects, so we actually shouldn't expect the ordinary arithmetic addition of velocities to hold in relativity; it's an approximation that's valid at velocities that are small compared to c .

A universal speed limit

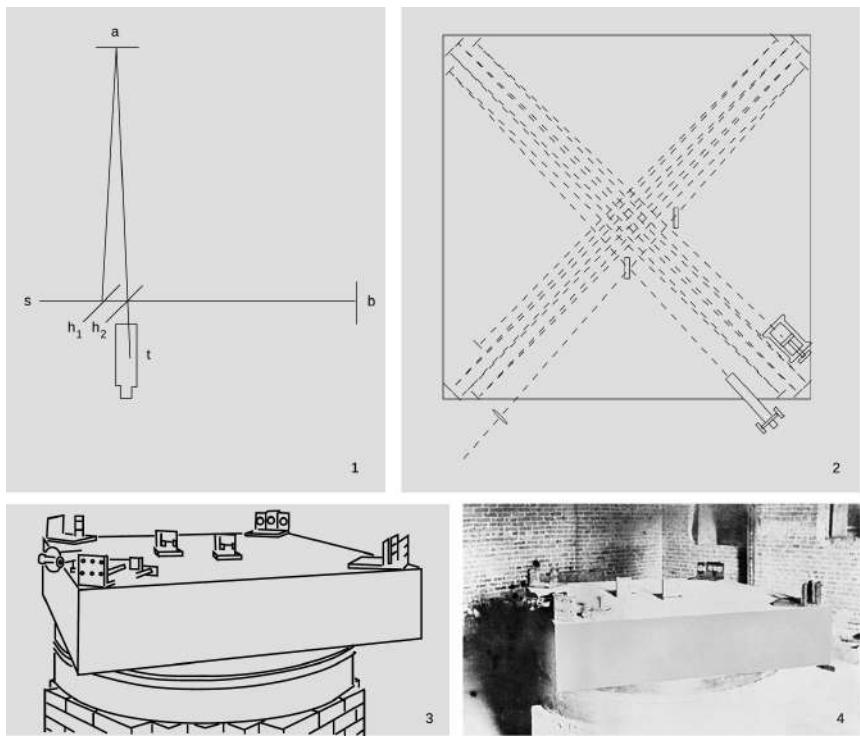
For example, suppose Janet takes a trip in a spaceship, and accelerates until she is moving at $0.6c$ relative to the earth. She then launches a space probe in the forward direction at a speed relative to her ship of $0.6c$. We might think that the probe was then moving at a velocity of $1.2c$, but in fact the answer is still less than c (problem 1, page 457). This is an example of a more general fact about relativity, which is that c represents a universal speed limit. This is required by causality, as shown in figure t.

Light travels at c .

Now consider a beam of light. We're used to talking casually about the "speed of light," but what does that really mean? Motion is relative, so normally if we want to talk about a velocity, we have to specify what it's measured relative to. A sound wave has a certain speed relative to the air, and a water wave has its own speed relative to the water. If we want to measure the speed of an ocean wave, for example, we should make sure to measure it in a frame of reference at rest relative to the water. But light isn't a vibration of a physical medium; it can propagate through the near-perfect vacuum of outer



t / A proof that causality imposes a universal speed limit. In the original frame of reference, represented by the square, event A happens a little before event B. In the new frame, shown by the parallelogram, A happens after $t = 0$, but B happens before $t = 0$; that is, B happens before A. The time ordering of the two events has been reversed. This can only happen because events A and B are very close together in time and fairly far apart in space. The line segment connecting A and B has a slope greater than 1, meaning that if we wanted to be present at both events, we would have to travel at a speed greater than c (which equals 1 in the units used on this graph). You will find that if you pick any two points for which the slope of the line segment connecting them is less than 1, you can never get them to straddle the new $t = 0$ line in this funny, time-reversed way. Since different observers disagree on the time order of events like A and B, causality requires that information never travel from A to B or from B to A; if it did, then we would have time-travel paradoxes. The conclusion is that c is the maximum speed of cause and effect in relativity.



u / The Michelson-Morley experiment, shown in photographs, and drawings from the original 1887 paper. 1. A simplified drawing of the apparatus. A beam of light from the source, s , is partially reflected and partially transmitted by the half-silvered mirror h_1 . The two half-intensity parts of the beam are reflected by the mirrors at a and b , reunited, and observed in the telescope, t . If the earth's surface was supposed to be moving through the ether, then the times taken by the two light waves to pass through the moving ether would be unequal, and the resulting time lag would be detectable by observing the interference between the waves when they were reunited. 2. In the real apparatus, the light beams were reflected multiple times. The effective length of each arm was increased to 11 meters, which greatly improved its sensitivity to the small expected difference in the speed of light. 3. In an earlier version of the experiment, they had run into problems with its "extreme sensitiveness to vibration," which was "so great that it was impossible to see the interference fringes except at brief intervals ... even at two o'clock in the morning." They therefore mounted the whole thing on a massive stone floating in a pool of mercury, which also made it possible to rotate it easily. 4. A photo of the apparatus.

space, as when rays of sunlight travel to earth. This seems like a paradox: light is supposed to have a specific speed, but there is no way to decide what frame of reference to measure it in. The way out of the paradox is that light must travel at a velocity equal to c . Since all observers agree on a velocity of c , regardless of their frame of reference, everything is consistent.

The Michelson-Morley experiment

The constancy of the speed of light had in fact already been observed when Einstein was an 8-year-old boy, but because nobody could figure out how to interpret it, the result was largely ignored. In 1887 Michelson and Morley set up a clever apparatus to measure any difference in the speed of light beams traveling east-west and north-south. The motion of the earth around the sun at 110,000 km/hour (about 0.01% of the speed of light) is to our west during the day. Michelson and Morley believed that light was a vibration of a mysterious medium called the ether, so they expected that the speed of light would be a fixed value relative to the ether. As the earth moved through the ether, they thought they would observe an effect on the velocity of light along an east-west line. For instance, if they released a beam of light in a westward direction during the day, they expected that it would move away from them at less than the normal speed because the earth was chasing it through the ether. They were surprised when they found that the expected 0.01% change in the speed of light did not occur.

The ring laser gyroscope

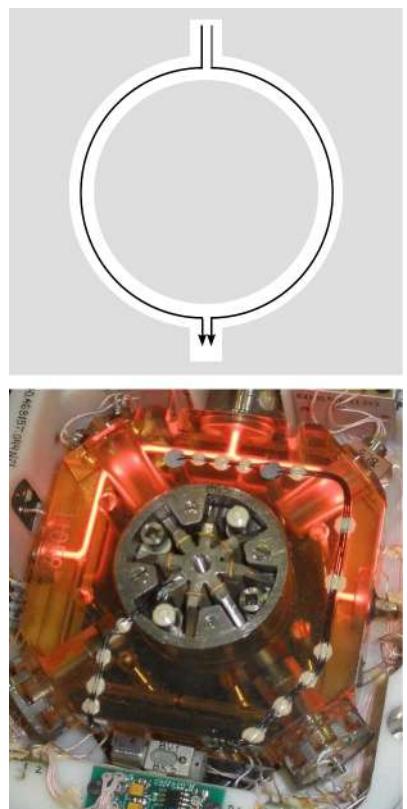
example 8

If you've flown in a jet plane, you can thank relativity for helping you to avoid crashing into a mountain or an ocean. Figure v shows a standard piece of navigational equipment called a ring laser gyroscope. A beam of light is split into two parts, sent around the perimeter of the device, and reunited. Since the speed of light is constant, we expect the two parts to come back together at the same time. If they don't, it's evidence that the device has been rotating. The plane's computer senses this and notes how much rotation has accumulated.

No frequency-dependence

example 9

Relativity has only one universal speed, so it requires that all light waves travel at the same speed, regardless of their frequency and wavelength. Presently the best experimental tests of the invariance of the speed of light with respect to wavelength come from astronomical observations of gamma-ray bursts, which are sudden outpourings of high-frequency light, believed to originate from a supernova explosion in another galaxy. One such observation, in 2009,³ found that the times of arrival of all the different frequencies in the burst differed by no more than 2 seconds out of a total time in flight on the order of ten billion years!



v / A ring laser gyroscope.

³<http://arxiv.org/abs/0908.1832>

Discussion Questions

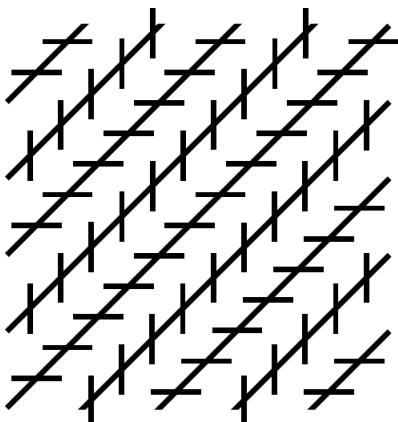
A A person in a spaceship moving at 99.9999999% of the speed of light relative to Earth shines a flashlight forward through dusty air, so the beam is visible. What does she see? What would it look like to an observer on Earth?

B A question that students often struggle with is whether time and space can really be distorted, or whether it just seems that way. Compare with optical illusions or magic tricks. How could you verify, for instance, that the lines in the figure are actually parallel? Are relativistic effects the same, or not?

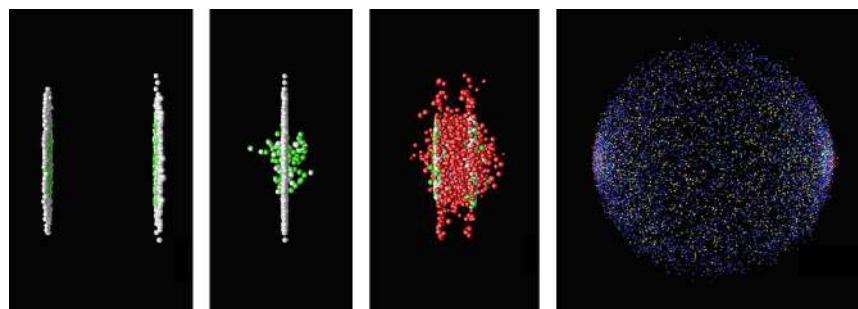
C On a spaceship moving at relativistic speeds, would a lecture seem even longer and more boring than normal?

D Mechanical clocks can be affected by motion. For example, it was a significant technological achievement to build a clock that could sail aboard a ship and still keep accurate time, allowing longitude to be determined. How is this similar to or different from relativistic time dilation?

E Figure r from page 409, depicting the collision of two nuclei at the RHIC accelerator, is reproduced below. What would the shapes of the two nuclei look like to a microscopic observer riding on the left-hand nucleus? To an observer riding on the right-hand one? Can they agree on what is happening? If not, why not — after all, shouldn't they see the same thing if they both compare the two nuclei side-by-side at the same instant in time?



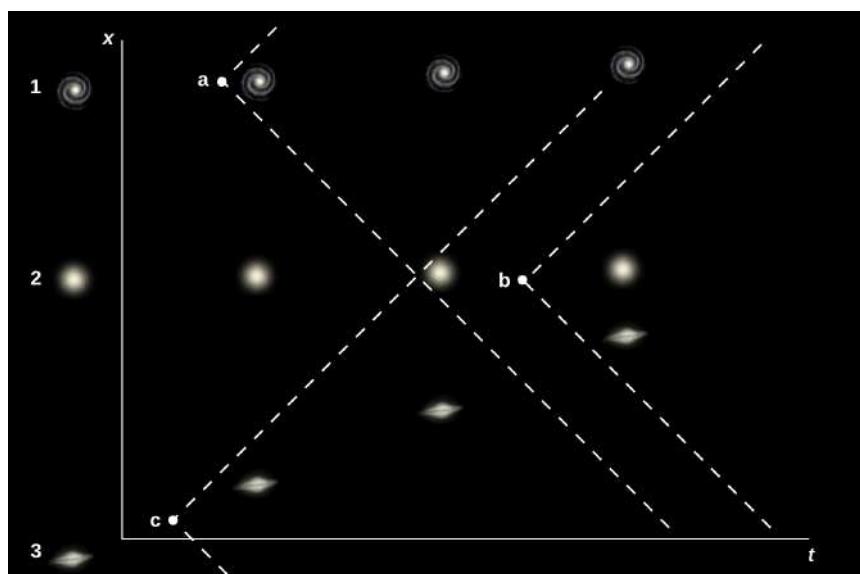
Discussion question B



w / Discussion question E: colliding nuclei show relativistic length contraction.

F If you stick a piece of foam rubber out the window of your car while driving down the freeway, the wind may compress it a little. Does it make sense to interpret the relativistic length contraction as a type of strain that pushes an object's atoms together like this? How does this relate to discussion question E?

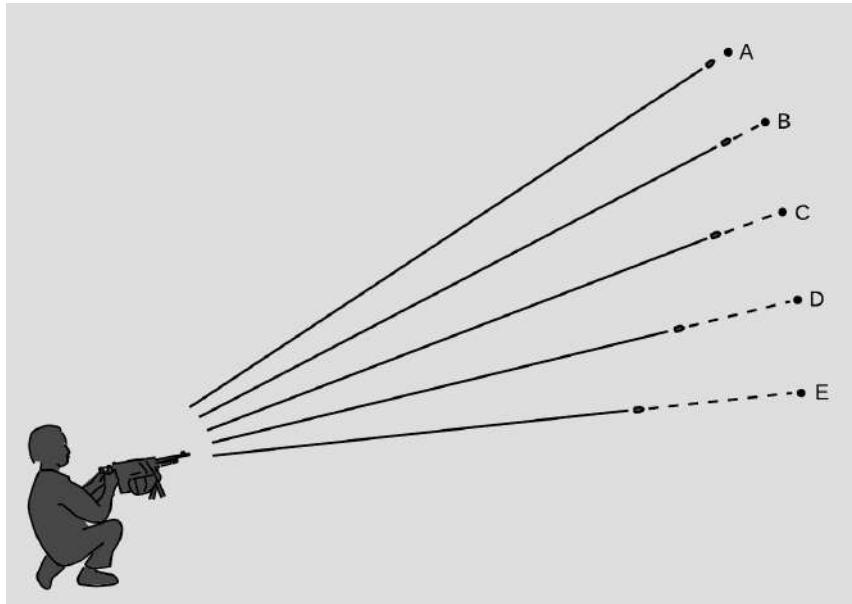
G The graph shows three galaxies. The axes are drawn according to an observer at rest relative to the galaxy 2, so that that galaxy is always at the same x coordinate. Intelligent species in the three different galaxies develop radio technology independently, and at some point each begins to actively send out signals in an attempt to communicate with other civilizations. Events a, b, and c mark the points at which these signals begin spreading out across the universe at the speed of light. Find the events at which the inhabitants of galaxy 2 detect the signals from galaxies 1 and 3. According to 2, who developed radio first, 1 or 3? On top of the graph, draw a new pair of position and time axes, for the frame in which galaxy 3 is at rest. According to 3, in what order did events a, b, and c happen?



Discussion question G.

H The machine-gunner in the figure sends out a spray of bullets. Suppose that the bullets are being shot into outer space, and that the distances traveled are trillions of miles (so that the human figure in the diagram is not to scale). After a long time, the bullets reach the points shown with dots which are all equally far from the gun. Their arrivals at those points are events A through E, which happen at different times. Sketch these events on a position-time graph. The chain of impacts extends across space at a speed greater than c . Does this violate special relativity?

Discussion question H.



7.2.4 No action at a distance

The Newtonian picture

The Newtonian picture of the universe has particles interacting with each other by exerting forces from a distance, and these forces are imagined to occur without any time delay. For example, suppose that super-powerful aliens, angered when they hear disco music in our AM radio transmissions, come to our solar system on a mission to cleanse the universe of our aesthetic contamination. They apply a force to our sun, causing it to go flying out of the solar system at a gazillion miles an hour. According to Newton's laws, the gravitational force of the sun on the earth will *immediately* start dropping off. This will be detectable on earth, and since sunlight takes eight minutes to get from the sun to the earth, the change in gravitational force will, according to Newton, be the first way in which earthlings learn the bad news — the sun will not visibly start receding until a little later. Although this scenario is fanciful, it shows a real feature of Newton's laws: that information can be transmitted from one place in the universe to another with zero time delay, so that transmission and reception occur at exactly the same instant. Newton was sharp enough to realize that this required a nontrivial assumption, which was that there was some completely objective and well-defined way of saying whether two things *happened* at exactly the same instant. He stated this assumption explicitly: “Absolute, true, and mathematical time, of itself, and from its own nature flows at a constant rate without regard to anything external...”

Time delays in forces exerted at a distance

Relativity forbids Newton's instantaneous action at a distance. For suppose that instantaneous action at a distance existed. It

would then be possible to send signals from one place in the universe to another without any time lag. This would allow perfect synchronization of all clocks. But the Hafele-Keating experiment demonstrates that clocks A and B that have been initially synchronized will drift out of sync if one is in motion relative to the other. With instantaneous transmission of signals, we could determine, without having to wait for A and B to be reunited, which was ahead and which was behind. Since they don't need to be reunited, neither one needs to undergo any acceleration; each clock can fix an inertial frame of reference, with a velocity vector that changes neither its direction nor its magnitude. But this violates the principle that constant-velocity motion is relative, because each clock can be considered to be at rest, in its own frame of reference. Since no experiment has ever detected any violation of the relativity of motion, we conclude that instantaneous action at a distance is impossible.

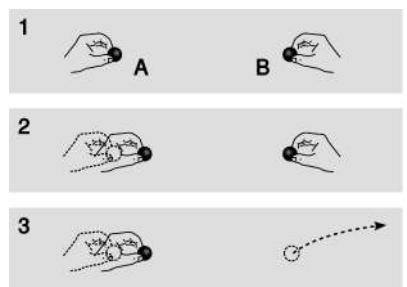
Since forces can't be transmitted instantaneously, it becomes natural to imagine force-effects spreading outward from their source like ripples on a pond, and we then have no choice but to impute some physical reality to these ripples. We call them fields, and they have their own independent existence. Gravity is transmitted through a field called the gravitational field. Besides gravity, there are other fundamental fields of force such as electricity and magnetism (). Ripples of the electric and magnetic fields turn out to be light waves. This tells us that the speed at which electric and magnetic field ripples spread must be c , and by an argument similar to the one in subsection 7.2.3 the same must hold for any other fundamental field, including the gravitational field.

Fields don't have to wiggle; they can hold still as well. The earth's magnetic field, for example, is nearly constant, which is why we can use it for direction-finding.

Even empty space, then, is not perfectly featureless. It has measurable properties. For example, we can drop a rock in order to measure the direction of the gravitational field, or use a magnetic compass to find the direction of the magnetic field. This concept made a deep impression on Einstein as a child. He recalled that when he was five years old, the gift of a magnetic compass convinced him that there was "something behind things, something deeply hidden."

More evidence that fields of force are real: they carry energy.

The smoking-gun argument for this strange notion of traveling force ripples comes from the fact that they carry energy. In figure z/1, Alice and Betty hold balls A and B at some distance from one another. These balls make a force on each other; it doesn't really matter for the sake of our argument whether this force is gravitational, electrical, or magnetic. Let's say it's electrical, i.e., that the balls have the kind of electrical *charge* that sometimes



z / Fields carry energy.

causes your socks to cling together when they come out of the clothes dryer. We'll say the force is repulsive, although again it doesn't really matter.

If Alice chooses to move her ball closer to Betty's, $z/2$, Alice will have to do some mechanical work against the electrical repulsion, burning off some of the calories from that chocolate cheesecake she had at lunch. This reduction in her body's chemical energy is offset by a corresponding increase in the electrical interaction energy. Not only that, but Alice feels the resistance stiffen as the balls get closer together and the repulsion strengthens. She has to do a little extra work, but this is all properly accounted for in the interaction energy.

But now suppose, $z/3$, that Betty decides to play a trick on Alice by tossing B far away just as Alice is getting ready to move A. We have already established that Alice can't feel B's motion instantaneously, so the electric forces must actually be propagated by an electric *field*. Of course this experiment is utterly impractical, but suppose for the sake of argument that the time it takes the change in the electric field to propagate across the diagram is long enough so that Alice can complete her motion before she feels the effect of B's disappearance. She is still getting stale information about B's position. As she moves A to the right, she feels a repulsion, because the field in her region of space is still the field caused by B in its *old* position. She has burned some chocolate cheesecake calories, and it appears that conservation of energy has been violated, because these calories can't be properly accounted for by any interaction with B, which is long gone.

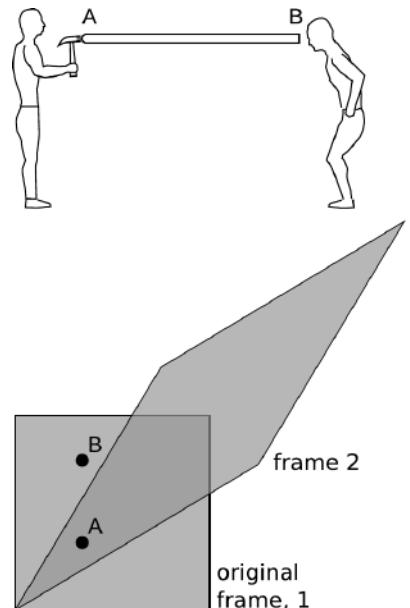
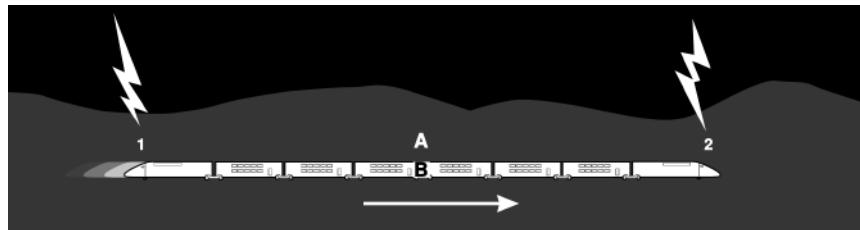
If we hope to preserve the law of conservation of energy, then the only possible conclusion is that the electric field itself carries away the cheesecake energy. In fact, this example represents an impractical method of transmitting radio waves. Alice does work on charge A, and that energy goes into the radio waves. Even if B had never existed, the radio waves would still have carried energy, and Alice would still have had to do work in order to create them.

Discussion Questions

A Amy and Bill are flying on spaceships in opposite directions at such high velocities that the relativistic effect on time's rate of flow is easily noticeable. Motion is relative, so Amy considers herself to be at rest and Bill to be in motion. She says that time is flowing normally for her, but Bill is slow. But Bill can say exactly the same thing. How can they *both* think the other is slow? Can they settle the disagreement by getting on the radio and seeing whose voice is normal and whose sounds slowed down and Darth-Vadery?



B The figure shows a famous thought experiment devised by Einstein. A train is moving at constant velocity to the right when bolts of lightning strike the ground near its front and back. Alice, standing on the dirt at the midpoint of the flashes, observes that the light from the two flashes arrives simultaneously, so she says the two strikes must have occurred simultaneously. Bob, meanwhile, is sitting aboard the train, at its middle. He passes by Alice at the moment when Alice later figures out that the flashes happened. Later, he receives flash 2, and then flash 1. He infers that since both flashes traveled half the length of the train, flash 2 must have occurred first. How can this be reconciled with Alice's belief that the flashes were simultaneous? Explain using a graph.



ab / Discussion question E.

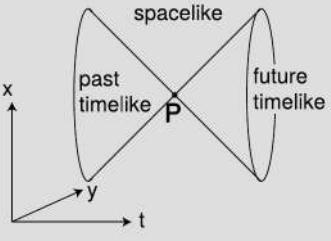
C Resolve the following paradox by drawing a spacetime diagram (i.e., a graph of x versus t). Andy and Beth are in motion relative to one another at a significant fraction of c . As they pass by each other, they exchange greetings, and Beth tells Andy that she is going to blow up a stick of dynamite one hour later. One hour later by Andy's clock, she still hasn't exploded the dynamite, and he says to himself, "She hasn't exploded it because of time dilation. It's only been 40 minutes for her." He now accelerates suddenly so that he's moving at the same velocity as Beth. The time dilation no longer exists. If he looks again, does he suddenly see the flash from the explosion? How can this be? Would he see her go through 20 minutes of her life in fast-motion?

D Use a graph to resolve the following relativity paradox. Relativity says that in one frame of reference, event A could happen before event B, but in someone else's frame B would come before A. How can this be? Obviously the two people could meet up at A and talk as they cruised past each other. Wouldn't they have to agree on whether B had already happened?

E The rod in the figure is perfectly rigid. At event A, the hammer strikes one end of the rod. At event B, the other end moves. Since the rod is perfectly rigid, it can't compress, so A and B are simultaneous. In frame 2, B happens before A. Did the motion at the right end cause the person on the left to decide to pick up the hammer and use it?

7.2.5 The light cone

Given an event P, we can now classify all the causal relationships in which P can participate. In Newtonian physics, these relationships fell into two classes: P could potentially cause any event that lay in its future, and could have been caused by any event in its past. In relativity, we have a three-way distinction rather than a two-way one. There is a third class of events that are too far away



ac / The light cone.

from P in space, and too close in time, to allow any cause and effect relationship, since causality's maximum velocity is c . Since we're working in units in which $c = 1$, the boundary of this set is formed by the lines with slope ± 1 on a (t, x) plot. This is referred to as the light cone, for reasons that become more visually obvious when we consider more than one spatial dimension, figure ac.

Events lying inside one another's light cones are said to have a timelike relationship. Events outside each other's light cones are spacelike in relation to one another, and in the case where they lie on the surfaces of each other's light cones the term is lightlike.

7.2.6 * The spacetime interval

The light cone is an object of central importance in both special and general relativity. It relates the *geometry* of spacetime to possible *cause-and-effect* relationships between events. This is fundamentally how relativity works: it's a geometrical theory of causality.

These ideas naturally lead us to ask what fruitful analogies we can form between the bizarre geometry of spacetime and the more familiar geometry of the Euclidean plane. The light cone cuts spacetime into different regions according to certain measurements of relationships between points (events). Similarly, a circle in Euclidean geometry cuts the plane into two parts, an interior and an exterior, according to the measurement of the distance from the circle's center. A circle stays the same when we rotate the plane. A light cone stays the same when we change frames of reference. Let's build up the analogy more explicitly.

Measurement in Euclidean geometry

We say that two line segments are congruent, $AB \cong CD$, if the distance between points A and B is the same as the distance between C and D, as measured by a rigid ruler.

Measurement in spacetime

We define $AB \cong CD$ if:

1. AB and CD are both spacelike, and the two distances are equal as measured by a rigid ruler, in a frame where the two events touch the ruler simultaneously.
2. AB and CD are both timelike, and the two time intervals are equal as measured by clocks moving inertially.
3. AB and CD are both lightlike.

The three parts of the relativistic version each require some justification.

Case 1 has to be the way it is because space is part of spacetime. In special relativity, this space is Euclidean, so the definition of congruence has to agree with the Euclidean definition, in the case

where it is possible to apply the Euclidean definition. The spacelike relation between the points is both necessary and sufficient to make this possible. If points A and B are spacelike in relation to one another, then a frame of reference exists in which they are simultaneous, so we can use a ruler that is at rest in that frame to measure their distance. If they are lightlike or timelike, then no such frame of reference exists. For example, there is no frame of reference in which Charles VII's restoration to the throne is simultaneous with Joan of Arc's execution, so we can't arrange for both of these events to touch the same ruler at the same time.

The definition in case 2 is the only sensible way to proceed if we are to respect the symmetric treatment of time and space in relativity. The timelike relation between the events is necessary and sufficient to make it possible for a clock to move from one to the other. It makes a difference that the clocks move inertially, because the twins in example 3 on p. 407 disagree on the clock time between the traveling twin's departure and return.

Case 3 may seem strange, since it says that *any* two lightlike intervals are congruent. But this is the only possible definition, because this case can be obtained as a limit of the timelike one. Suppose that AB is a timelike interval, but in the planet earth's frame of reference it would be necessary to travel at almost the speed of light in order to reach B from A. The required speed is less than c (i.e., less than 1) by some tiny amount ϵ . In the earth's frame, the clock referred to in the definition suffers extreme time dilation. The time elapsed on the clock is very small. As ϵ approaches zero, and the relationship between A and B approaches a lightlike one, this clock time approaches zero. In this sense, the relativistic notion of "distance" is very different from the Euclidean one. In Euclidean geometry, the distance between two points can only be zero if they are the same point.

The case splitting involved in the relativistic definition is a little ugly. Having worked out the physical interpretation, we can now consolidate the definition in a nicer way by appealing to Cartesian coordinates.

Cartesian definition of distance in Euclidean geometry

Given a vector $(\Delta x, \Delta y)$ from point A to point B, the square of the distance between them is defined as $\overline{AB}^2 = \Delta x^2 + \Delta y^2$.

Definition of the interval in relativity

Given points separated by coordinate differences Δx , Δy , Δz , and Δt , the spacetime interval \mathcal{I} (cursive letter "I") between them is defined as $\mathcal{I} = \Delta t^2 - \Delta x^2 - \Delta y^2 - \Delta z^2$.

This is stated in natural units, so all four terms on the right-hand side have the same units; in metric units with $c \neq 1$, appropriate factors of c should be inserted in order to make the units of the

terms agree. The interval \mathcal{I} is positive if AB is timelike (regardless of which event comes first), zero if lightlike, and negative if spacelike. Since \mathcal{I} can be negative, we can't in general take its square root and define a real number \overline{AB} as in the Euclidean case. When the interval is timelike, we can interpret $\sqrt{\mathcal{I}}$ as a time, and when it's spacelike we can take $\sqrt{-\mathcal{I}}$ to be a distance.

The Euclidean definition of distance (i.e., the Pythagorean theorem) is useful because it gives the same answer regardless of how we rotate the plane. Although it is stated in terms of a certain coordinate system, its result is unambiguously defined because it is the same regardless of what coordinate system we arbitrarily pick. Similarly, \mathcal{I} is useful because, as proved in example 11 below, it is the same regardless of our frame of reference, i.e., regardless of our choice of coordinates.

Pioneer 10

example 10

▷ The Pioneer 10 space probe was launched in 1972, and in 1973 was the first craft to fly by the planet Jupiter. It crossed the orbit of the planet Neptune in 1983, after which telemetry data were received until 2002. The following table gives the spacecraft's position relative to the sun at exactly midnight on January 1, 1983 and January 1, 1995. The 1983 date is taken to be $t = 0$.

t (s)	x	y	z
0	1.784×10^{12} m	3.951×10^{12} m	0.237×10^{12} m
3.7869120000×10^8 s	2.420×10^{12} m	8.827×10^{12} m	0.488×10^{12} m

Compare the time elapsed on the spacecraft to the time in a frame of reference tied to the sun.

▷ We can convert these data into natural units, with the distance unit being the second (i.e., a light-second, the distance light travels in one second) and the time unit being seconds. Converting and carrying out this subtraction, we have:

Δt (s)	Δx	Δy	Δz
3.7869120000×10^8 s	0.2121×10^4 s	1.626×10^4 s	0.084×10^4 s

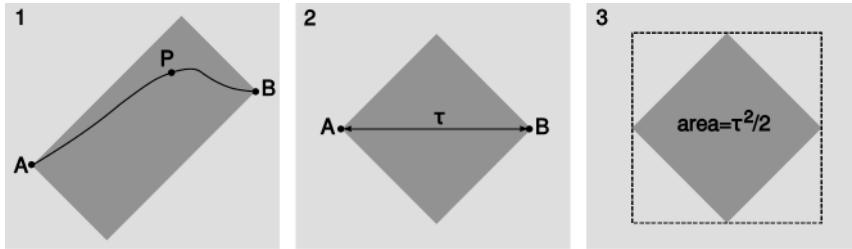
Comparing the exponents of the temporal and spatial numbers, we can see that the spacecraft was moving at a velocity on the order of 10^{-4} of the speed of light, so relativistic effects should be small but not completely negligible.

Since the interval is timelike, we can take its square root and interpret it as the time elapsed on the spacecraft. The result is $\sqrt{\mathcal{I}} = 3.786911996 \times 10^8$ s. This is 0.4 s less than the time elapsed in the sun's frame of reference.

Invariance of the interval

example 11

In this example we prove that the interval is the same regard-



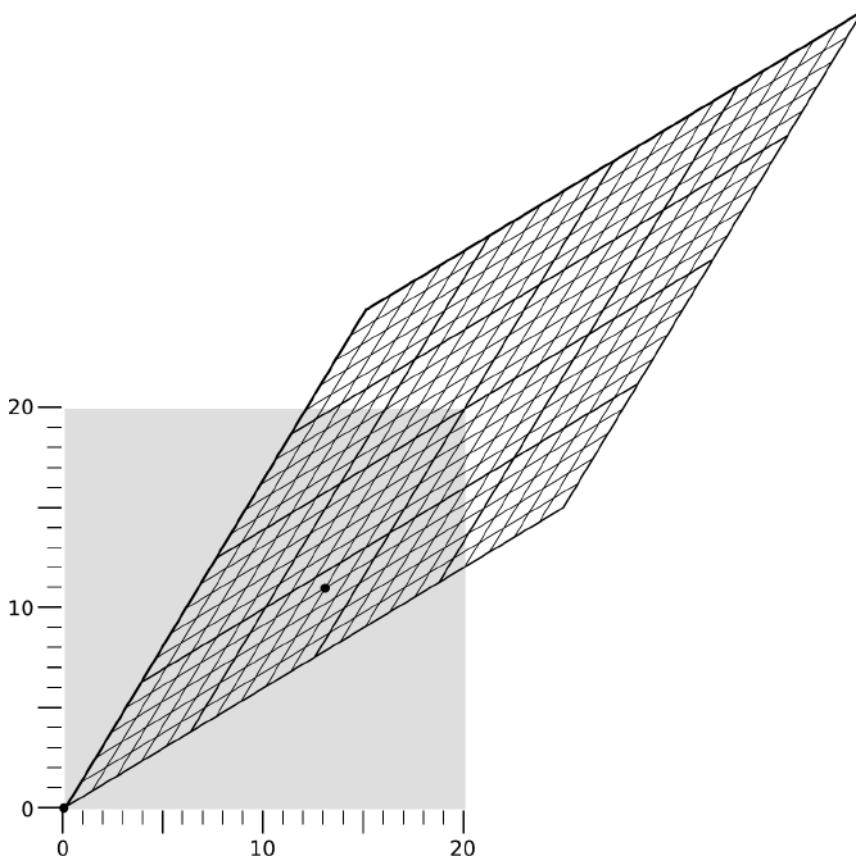
less of what frame of reference we compute it in. This is called “Lorentz invariance.” The proof is limited to the timelike case. Given events A and B, construct the light-rectangle as defined in figure ad/1. On p. 405 we proved that the Lorentz transformation doesn’t change the area of a shape in the x - t plane. Therefore the area of this rectangle is unchanged if we switch to the frame of reference ad/2, in which A and B occurred at the same location and were separated by a time interval τ . This area equals half the interval \mathcal{I} between A and B. But a straightforward calculation shows that the rectangle in ad/1 also has an area equal to half the interval calculated in *that* frame. Since the area in any frame equals half the interval, and the area is the same in all frames, the interval is equal in all frames as well.

ad / Light-rectangles, example 11.

1. The gray light-rectangle represents the set of all events such as P that could be visited after A and before B.

2. The rectangle becomes a square in the frame in which A and B occur at the same location in space.

3. The area of the dashed square is τ^2 , so the area of the gray square is $\tau^2/2$.



ae / Example 12.

A numerical example of invariance *example 12*

Figure ae shows two frames of reference in motion relative to one another at $v = 3/5$. (For this velocity, the stretching and squishing of the main diagonals are both by a factor of 2.) Events are marked at coordinates that in the frame represented by the square are

$$(t, x) = (0, 0) \quad \text{and} \\ (t, x) = (13, 11).$$

The interval between these events is $13^2 - 11^2 = 48$. In the frame represented by the parallelogram, the same two events lie at coordinates

$$(t', x') = (0, 0) \quad \text{and} \\ (t', x') = (8, 4).$$

Calculating the interval using these values, the result is $8^2 - 4^2 = 48$, which comes out the same as in the other frame.

7.2.7 * Four-vectors and the inner product

Example 10 makes it natural that we define a type of vector with four components, the first one relating to time and the others being spatial. These are known as four-vectors. It's clear how we should define the equivalent of a dot product in relativity:

$$\mathbf{A} \cdot \mathbf{B} = A_t B_t - A_x B_x - A_y B_y - A_z B_z$$

The term “dot product” has connotations of referring only to three-vectors, so the operation of taking the scalar product of two four-vectors is usually referred to instead as the “inner product.” The spacetime interval can then be thought of as the inner product of a four-vector with itself. We care about the relativistic inner product for exactly the same reason we care about its Euclidean version; both are scalars, so they have a fixed value regardless of what coordinate system we choose.

The twin paradox

example 13

Alice and Betty are identical twins. Betty goes on a space voyage at relativistic speeds, traveling away from the earth and then turning around and coming back. Meanwhile, Alice stays on earth. When Betty returns, she is younger than Alice because of relativistic time dilation (example 3, p. 407).

But isn't it valid to say that Betty's spaceship is standing still and the earth moving? In that description, wouldn't Alice end up younger and Betty older? This is referred to as the “twin paradox.” It can't really be a paradox, since it's exactly what was observed in the Hafele-Keating experiment (p. 397).

Betty's track in the x - t plane (her “world-line” in relativistic jargon) consists of vectors **b** and **c** strung end-to-end (figure af). We could adopt a frame of reference in which Betty was at rest during **b** (i.e., $b_x = 0$), but there is no frame in which **b** and **c** are parallel, so there is no frame in which Betty was at rest during *both* **b** and **c**. This resolves the paradox.

We have already established by other methods that Betty ages less than Alice, but let's see how this plays out in a simple numerical example. Omitting units and making up simple numbers, let's say that the vectors in figure af are

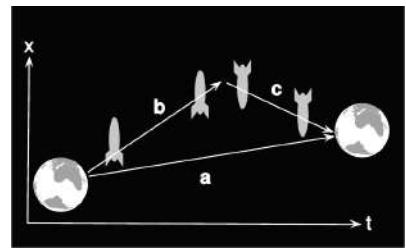
$$\mathbf{a} = (6, 1)$$

$$\mathbf{b} = (3, 2)$$

$$\mathbf{c} = (3, -1),$$

where the components are given in the order (t, x) . The time experienced by Alice is then

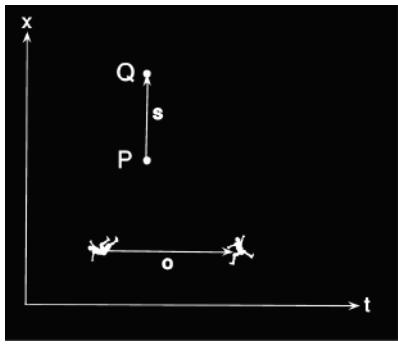
$$|\mathbf{a}| = \sqrt{6^2 - 1^2} = 5.9,$$



af / Example 13.

which is greater than the Betty's elapsed time

$$|\mathbf{b}| + |\mathbf{c}| = \sqrt{3^2 - 2^2} + \sqrt{3^2 - (-1)^2} = 5.1.$$



ag / Example 14.

Simultaneity using inner products

example 14

Suppose that an observer O moves inertially along a vector \mathbf{o} , and let the vector separating two events P and Q be \mathbf{s} . O judges these events to be simultaneous if $\mathbf{o} \cdot \mathbf{s} = 0$. To see why this is true, suppose we pick a coordinate system as defined by O. In this coordinate system, O considers herself to be at rest, so she says her vector has only a time component, $\mathbf{o} = (\Delta t, 0, 0, 0)$. If she considers P and Q to be simultaneous, then the vector from P to Q is of the form $(0, \Delta x, \Delta y, \Delta z)$. The inner product is then zero, since each of the four terms vanishes. Since the inner product is independent of the choice of coordinate system, it doesn't matter that we chose one tied to O herself. Any other observer O' can look at O's motion, note that $\mathbf{o} \cdot \mathbf{s} = 0$, and infer that O must consider P and Q to be simultaneous, even if O' says they weren't.

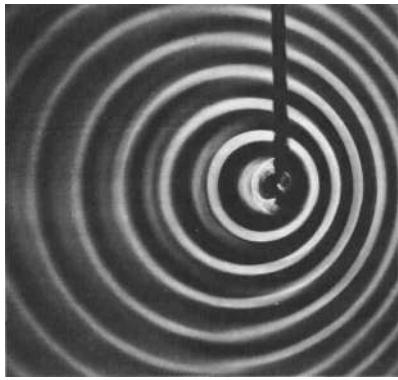
7.2.8 * Doppler shifts of light and addition of velocities

When Doppler shifts happen to ripples on a pond or the sound waves from an airplane, they can depend on the relative motion of three different objects: the source, the receiver, and the medium. But light waves don't have a medium. Therefore Doppler shifts of light can only depend on the relative motion of the source and observer.

One simple case is the one in which the relative motion of the source and the receiver is perpendicular to the line connecting them. That is, the motion is transverse. Nonrelativistic Doppler shifts happen because the distance between the source and receiver is changing, so in nonrelativistic physics we don't expect any Doppler shift at all when the motion is transverse, and this is what is in fact observed to high precision. For example, the photo shows shortened and lengthened wavelengths to the right and left, along the source's line of motion, but an observer above or below the source measures just the normal, unshifted wavelength and frequency. But relativistically, we have a time dilation effect, so for light waves emitted transversely, there is a Doppler shift of $1/\gamma$ in frequency (or γ in wavelength).

The other simple case is the one in which the relative motion of the source and receiver is longitudinal, i.e., they are either approaching or receding from one another. For example, distant galaxies are receding from our galaxy due to the expansion of the universe, and this expansion was originally detected because Doppler shifts toward the red (low-frequency) end of the spectrum were observed.

Nonrelativistically, we would expect the light from such a galaxy to be Doppler shifted down in frequency by some factor, which



ah / The pattern of waves made by a point source moving to the right across the water. Note the shorter wavelength of the forward-emitted waves and the longer wavelength of the backward-going ones.

would depend on the relative velocities of three different objects: the source, the wave's medium, and the receiver. Relativistically, things get simpler, because light isn't a vibration of a physical medium, so the Doppler shift can only depend on a single velocity v , which is the rate at which the separation between the source and the receiver is increasing.

The square in figure aj is the “graph paper” used by someone who considers the source to be at rest, while the parallelogram plays a similar role for the receiver. The figure is drawn for the case where $v = 3/5$ (in units where $c = 1$), and in this case the stretch factor of the long diagonal is 2. To keep the area the same, the short diagonal has to be squished to half its original size. But now it's a matter of simple geometry to show that OP equals half the width of the square, and this tells us that the Doppler shift is a factor of $1/2$ in frequency. That is, the squish factor of the short diagonal is interpreted as the Doppler shift. To get this as a general equation for velocities other than $3/5$, one can show by straightforward fiddling with the result of part c of problem 7 on p. 458 that the Doppler shift is

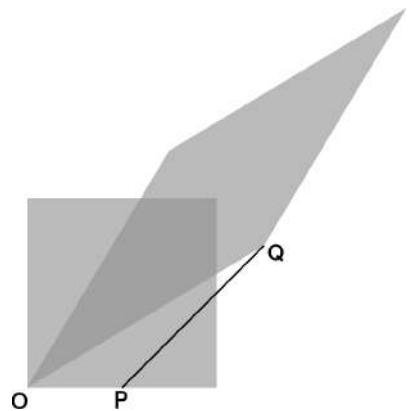
$$D(v) = \sqrt{\frac{1-v}{1+v}}.$$

Here $v > 0$ is the case where the source and receiver are getting farther apart, $v < 0$ the case where they are approaching. (This is the opposite of the sign convention used in subsection 6.1.5. It is convenient to change conventions here so that we can use positive values of v in the case of cosmological red-shifts, which are the most important application.)

Suppose that Alice stays at home on earth while her twin Betty takes off in her rocket ship at $3/5$ of the speed of light. When I first learned relativity, the thing that caused me the most pain was understanding how each observer could say that the other was the one whose time was slow. It seemed to me that if I could take a pill that would speed up my mind and my body, then naturally I would see everybody *else* as being *slow*. Shouldn't the same apply to relativity? But suppose Alice and Betty get on the radio and try to settle who is the fast one and who is the slow one. Each twin's voice sounds sloooowwww to the other. If Alice claps her hands twice, at a time interval of one second by her clock, Betty hears the hand-claps coming over the radio two seconds apart, but the situation is exactly symmetric, and Alice hears the same thing if Betty claps. Each twin analyzes the situation using a diagram identical to aj, and attributes her sister's observations to a complicated combination of time distortion, the time taken by the radio signals to propagate, and the motion of her twin relative to her.



ai / A graphical representation of the Lorentz transformation for a velocity of $(3/5)c$. The long diagonal is stretched by a factor of two, the short one is half its former length, and the area is the same as before.



aj / At event O, the source and the receiver are on top of each other, so as the source emits a wave crest, it is received without any time delay. At P, the source emits another wave crest, and at Q the receiver receives it.

self-check B

Turn your book upside-down and reinterpret figure aj. ▷ Answer, p. 1062

A symmetry property of the Doppler effect example 15

Suppose that A and B are at rest relative to one another, but C is moving along the line between A and B. A transmits a signal to C, who then retransmits it to B. The signal accumulates two Doppler shifts, and the result is their product $D(v)D(-v)$. But this product must equal 1, so we must have $D(-v)D(v) = 1$, which can be verified directly from the equation.

The Ives-Stilwell experiment example 16

The result of example 15 was the basis of one of the earliest laboratory tests of special relativity, by Ives and Stilwell in 1938. They observed the light emitted by excited by a beam of H_2^+ and H_3^+ ions with speeds of a few tenths of a percent of c . Measuring the light from both ahead of and behind the beams, they found that the product of the Doppler shifts $D(v)D(-v)$ was equal to 1, as predicted by relativity. If relativity had been false, then one would have expected the product to differ from 1 by an amount that would have been detectable in their experiment. In 2003, Saathoff et al. carried out an extremely precise version of the Ives-Stilwell technique with Li^+ ions moving at 6.4% of c . The frequencies observed, in units of MHz, were:

$$\begin{aligned}f_0 &= 546466918.8 \pm 0.4 \\&\quad (\text{unshifted frequency}) \\f_0 D(-v) &= 582490203.44 \pm .09 \\&\quad (\text{shifted frequency, forward}) \\f_0 D(v) &= 512671442.9 \pm 0.5 \\&\quad (\text{shifted frequency, backward}) \\ \sqrt{f_0 D(-v) \cdot f_0 D(v)} &= 546466918.6 \pm 0.3\end{aligned}$$

The results show incredibly precise agreement between f_0 and $\sqrt{f_0 D(-v) \cdot f_0 D(v)}$, as expected relativistically because $D(v)D(-v)$ is supposed to equal 1. The agreement extends to 9 significant figures, whereas if relativity had been false there should have been a relative disagreement of about $v^2 = .004$, i.e., a discrepancy in the third significant figure. The spectacular agreement with theory has made this experiment a lightning rod for anti-relativity kooks.

We saw on p. 411 that relativistic velocities should not be expected to be exactly additive, and problem 1 on p. 457 verifies this in the special case where A moves relative to B at $0.6c$ and B relative to C at $0.6c$ — the result *not* being $1.2c$. The relativistic Doppler shift provides a simple way of deriving a general equation for the relativistic combination of velocities; problem 17 on p. 461 guides you through the steps of this derivation, and the result is given on p. 1046.

7.3 Dynamics

So far we have said nothing about how to predict motion in relativity. Do Newton's laws still work? Do conservation laws still apply? The answer is yes, but many of the definitions need to be modified, and certain entirely new phenomena occur, such as the equivalence of energy and mass, as described by the famous equation $E = mc^2$.

7.3.1 Momentum

Consider the following scheme for traveling faster than the speed of light. The basic idea can be demonstrated by dropping a ping-pong ball and a baseball stacked on top of each other like a snowman. They separate slightly in mid-air, and the baseball therefore has time to hit the floor and rebound before it collides with the ping-pong ball, which is still on the way down. The result is a surprise if you haven't seen it before: the ping-pong ball flies off at high speed and hits the ceiling! A similar fact is known to people who investigate the scenes of accidents involving pedestrians. If a car moving at 90 kilometers per hour hits a pedestrian, the pedestrian flies off at nearly double that speed, 180 kilometers per hour. Now suppose the car was moving at 90 percent of the speed of light. Would the pedestrian fly off at 180% of c ?

To see why not, we have to back up a little and think about where this speed-doubling result comes from. For any collision, there is a special frame of reference, the center-of-mass frame, in which the two colliding objects approach each other, collide, and rebound with their velocities reversed. In the center-of-mass frame, the total momentum of the objects is zero both before and after the collision.

a / An unequal collision, viewed in the center-of-mass frame, 1, and in the frame where the small ball is initially at rest, 2. The motion is shown as it would appear on the film of an old-fashioned movie camera, with an equal amount of time separating each frame from the next. Film 1 was made by a camera that tracked the center of mass, film 2 by one that was initially tracking the small ball, and kept on moving at the same speed after the collision.

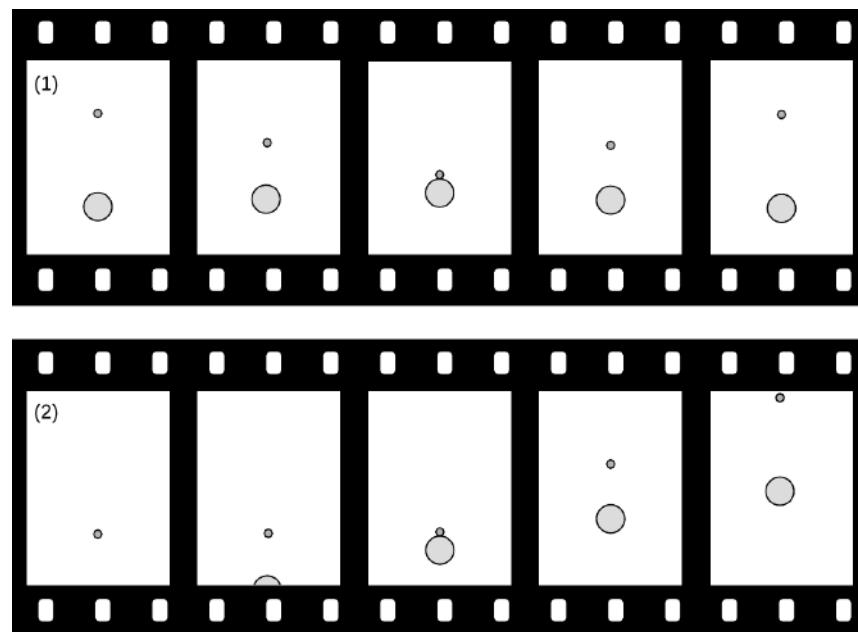


Figure a/1 shows such a frame of reference for objects of very unequal mass. Before the collision, the large ball is moving relatively slowly toward the top of the page, but because of its greater mass, its momentum cancels the momentum of the smaller ball, which is moving rapidly in the opposite direction. The total momentum is zero. After the collision, the two balls just reverse their directions of motion. We know that this is the right result for the outcome of the collision because it conserves both momentum and kinetic energy, and everything not forbidden is compulsory, i.e., in any experiment, there is only one possible outcome, which is the one that obeys all the conservation laws.

self-check C

How do we know that momentum and kinetic energy are conserved in figure a/1? ▷ Answer, p. 1062

Let's make up some numbers as an example. Say the small ball has a mass of 1 kg, the big one 8 kg. In frame 1, let's make the velocities as follows:

	before the collision	after the collision
•	-0.8	0.8
○	0.1	-0.1

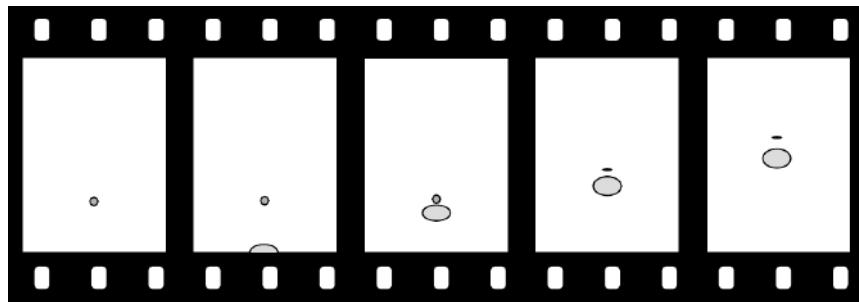
Figure a/2 shows the same collision in a frame of reference where the small ball was initially at rest. To find all the velocities in this frame, we just add 0.8 to all the ones in the previous table.

	before the collision	after the collision
•	0	1.6
○	0.9	0.7

In this frame, as expected, the small ball flies off with a velocity, 1.6, that is almost twice the initial velocity of the big ball, 0.9.

If all those velocities were in meters per second, then that's exactly what happened. But what if all these velocities were in units of the speed of light? Now it's no longer a good approximation just to add velocities. We need to combine them according to the relativistic rules. For instance, the technique used in problem 1 on p. 457 can be used to show that combining a velocity of 0.8 times the speed of light with another velocity of 0.8 results in 0.98, not 1.6. The results are very different:

	before the collision	after the collision
•	0	0.98
○	0.83	0.76

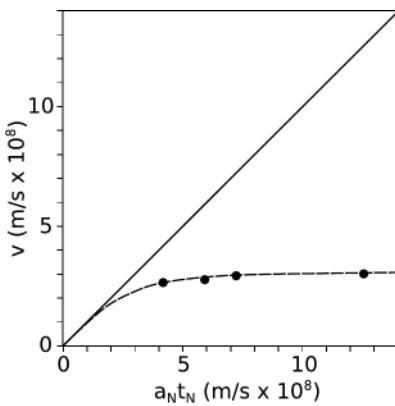


b / An 8-kg ball moving at 83% of the speed of light hits a 1-kg ball. The balls appear foreshortened due to the relativistic distortion of space.

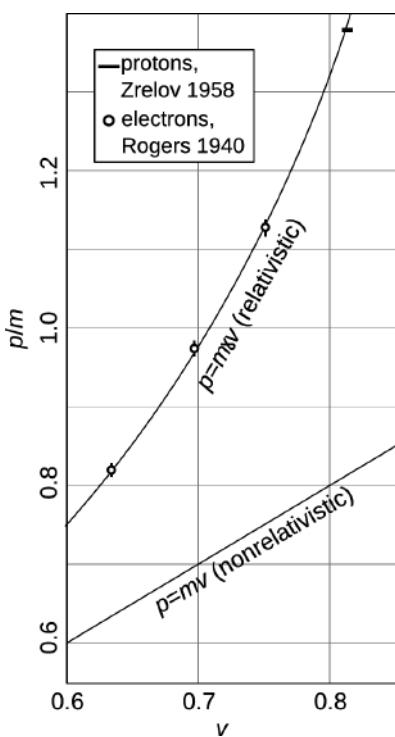
We can interpret this as follows. Figure a/1 is one in which the big ball is moving fairly slowly. This is very nearly the way the scene would be seen by an ant standing on the big ball. According to an observer in frame b, however, both balls are moving at nearly the speed of light after the collision. Because of this, the balls appear foreshortened, but the distance between the two balls is also shortened. To this observer, it seems that the small ball isn't pulling away from the big ball very fast.

Now here's what's interesting about all this. The outcome shown in figure a/2 was supposed to be the only one possible, the only one that satisfied both conservation of energy and conservation of momentum. So how can the *different* result shown in figure b be possible? The answer is that relativistically, momentum must not equal mv . The old, familiar definition is only an approximation that's valid at low speeds. If we observe the behavior of the small ball in figure b, it looks as though it somehow had some extra inertia. It's as though a football player tried to knock another player down without realizing that the other guy had a three-hundred-pound bag full of lead shot hidden under his uniform — he just doesn't seem to react to the collision as much as he should. As proved in section 7.3.4, this extra inertia is described by redefining momentum as

$$p = m\gamma v.$$



c / Example 17.



d / Two early high-precision tests of the relativistic equation $p = m\gamma v$ for momentum. Graphing p/m rather than p allows the data for electrons and protons to be placed on the same graph. Natural units are used, so that the horizontal axis is the velocity in units of c , and the vertical axis is the unitless quantity p/mc . The very small error bars for the data point from Zrelov are represented by the height of the black rectangle.

At very low velocities, γ is close to 1, and the result is very nearly mv , as demanded by the correspondence principle. But at very high velocities, γ gets very big — the small ball in figure b has a γ of 5.0, and therefore has five times more inertia than we would expect nonrelativistically.

This also explains the answer to another paradox often posed by beginners at relativity. Suppose you keep on applying a steady force to an object that's already moving at $0.9999c$. Why doesn't it just keep on speeding up past c ? The answer is that force is the rate of change of momentum. At $0.9999c$, an object already has a γ of 71, and therefore has already sucked up 71 times the momentum you'd expect at that speed. As its velocity gets closer and closer to c , its γ approaches infinity. To move at c , it would need an infinite momentum, which could only be caused by an infinite force.

Push as hard as you like ...

example 17

We don't have to depend on our imaginations to see what would happen if we kept on applying a force to an object indefinitely and tried to accelerate it past c . A nice experiment of this type was done by Bertozzi in 1964. In this experiment, electrons were accelerated by an electric field E through a distance ℓ_1 . Applying Newton's laws gives Newtonian predictions a_N for the acceleration and t_N for the time required.⁴

The electrons were then allowed to fly down a pipe for a further distance $\ell_2 = 8.4$ m without being acted on by any force. The time of flight t_2 for this second distance was used to find the final velocity $v = \ell_2/t_2$ to which they had actually been accelerated.

Figure c shows the results.⁵ According to Newton, an acceleration a_N acting for a time t_N should produce a final velocity $a_N t_N$. The solid line in the graph shows the prediction of Newton's laws, which is that a constant force exerted steadily over time will produce a velocity that rises linearly and without limit.

The experimental data, shown as black dots, clearly tell a different story. The velocity never goes above a certain maximum value, which we identify as c . The dashed line shows the predictions of special relativity, which are in good agreement with the experimental results.

Figure d shows experimental data confirming the relativistic equation for momentum.

⁴Newton's second law gives $a_N = F/m = eE/m$. The constant-acceleration equation $\Delta x = (1/2)at^2$ then gives $t_N = \sqrt{2m\ell_1/eE}$.

⁵To make the low-energy portion of the graph legible, Bertozzi's highest-energy data point is omitted.

7.3.2 Equivalence of mass and energy

Now we're ready to see why mass and energy must be equivalent as claimed in the famous $E = mc^2$. So far we've only considered collisions in which none of the kinetic energy is converted into any other form of energy, such as heat or sound. Let's consider what happens if a blob of putty moving at velocity v hits another blob that is initially at rest, sticking to it. The nonrelativistic result is that to obey conservation of momentum the two blobs must fly off together at $v/2$. Half of the initial kinetic energy has been converted to heat.⁶

Relativistically, however, an interesting thing happens. A hot object has more momentum than a cold object! This is because the relativistically correct expression for momentum is $m\gamma v$, and the more rapidly moving atoms in the hot object have higher values of γ . In our collision, the final combined blob must therefore be moving a little more slowly than the expected $v/2$, since otherwise the final momentum would have been a little greater than the initial momentum. To an observer who believes in conservation of momentum and knows only about the overall motion of the objects and not about their heat content, the low velocity after the collision would seem to be the result of a magical change in the mass, as if the mass of two combined, hot blobs of putty was more than the sum of their individual masses.

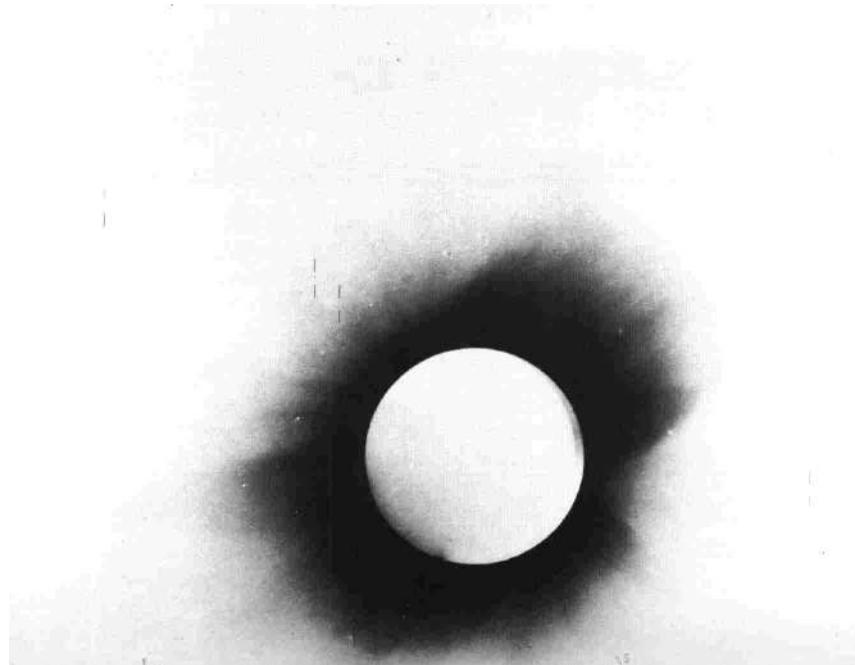
Now we know that the masses of all the atoms in the blobs must be the same as they always were. The change is due to the change in γ with heating, not to a change in mass. The heat energy, however, seems to be acting as if it was equivalent to some extra mass.

But this whole argument was based on the fact that heat is a form of kinetic energy at the atomic level. Would $E = mc^2$ apply to other forms of energy as well? Suppose a rocket ship contains some electrical energy stored in a battery. If we believed that $E = mc^2$ applied to forms of kinetic energy but not to electrical energy, then we would have to believe that the pilot of the rocket could slow the ship down by using the battery to run a heater! This would not only be strange, but it would violate the principle of relativity, because the result of the experiment would be different depending on whether the ship was at rest or not. The only logical conclusion is that all forms of energy are equivalent to mass. Running the heater then has no effect on the motion of the ship, because the total energy in the ship was unchanged; one form of energy (electrical) was simply converted to another (heat).

The equation $E = mc^2$ tells us how much energy is equivalent

⁶A double-mass object moving at half the speed does not have the same kinetic energy. Kinetic energy depends on the square of the velocity, so cutting the velocity in half reduces the energy by a factor of $1/4$, which, multiplied by the doubled mass, makes $1/2$ the original energy.

to how much mass: the conversion factor is the square of the speed of light, c . Since c a big number, you get a really really big number when you multiply it by itself to get c^2 . This means that even a small amount of mass is equivalent to a very large amount of energy.



LIGHTS ALL ASKEW IN THE HEAVENS

**Men of Science More or Less
Agog Over Results of Eclipse
Observations.**

EINSTEIN THEORY TRIUMPHS

**Stars Not Where They Seemed
or Were Calculated to be,
but Nobody Need Worry.**

A BOOK FOR 12 WISE MEN

**No More in All the World Could
Comprehend It, Said Einstein When
His Daring Publishers Accepted It.**

f / A New York Times headline from November 10, 1919, describing the observations discussed in example 18.

e / Example 18, page 434.

Gravity bending light

example 18

Gravity is a universal attraction between things that have mass, and since the energy in a beam of light is equivalent to some very small amount of mass, we expect that light will be affected by gravity, although the effect should be very small. The first important experimental confirmation of relativity came in 1919 when stars next to the sun during a solar eclipse were observed to have shifted a little from their ordinary position. (If there was no eclipse, the glare of the sun would prevent the stars from being observed.) Starlight had been deflected by the sun's gravity. Figure e is a photographic negative, so the circle that appears bright is actually the dark face of the moon, and the dark area is really the bright corona of the sun. The stars, marked by lines above and below them, appeared at positions slightly different than their normal ones.

Black holes

example 19

A star with sufficiently strong gravity can prevent light from leaving. Quite a few black holes have been detected via their gravitational forces on neighboring stars or clouds of gas and dust.

You've learned about conservation of mass and conservation of energy, but now we see that they're not even separate conservation laws. As a consequence of the theory of relativity, mass and energy are equivalent, and are not separately conserved — one can be converted into the other. Imagine that a magician waves his wand, and changes a bowl of dirt into a bowl of lettuce. You'd be impressed, because you were expecting that both dirt and lettuce would be conserved quantities. Neither one can be made to vanish, or to appear out of thin air. However, there are processes that can change one into the other. A farmer changes dirt into lettuce, and a compost heap changes lettuce into dirt. At the most fundamental level, lettuce and dirt aren't really different things at all; they're just collections of the same kinds of atoms — carbon, hydrogen, and so on. Because mass and energy are like two different sides of the same coin, we may speak of mass-energy, a single conserved quantity, found by adding up all the mass and energy, with the appropriate conversion factor: $E + mc^2$.

A rusting nail

example 20

- ▷ An iron nail is left in a cup of water until it turns entirely to rust. The energy released is about 0.5 MJ. In theory, would a sufficiently precise scale register a change in mass? If so, how much?
- ▷ The energy will appear as heat, which will be lost to the environment. The total mass-energy of the cup, water, and iron will indeed be lessened by 0.5 MJ. (If it had been perfectly insulated, there would have been no change, since the heat energy would have been trapped in the cup.) The speed of light is $c = 3 \times 10^8$ meters per second, so converting to mass units, we have

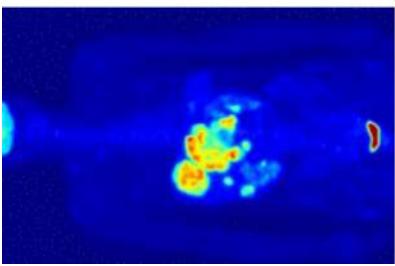
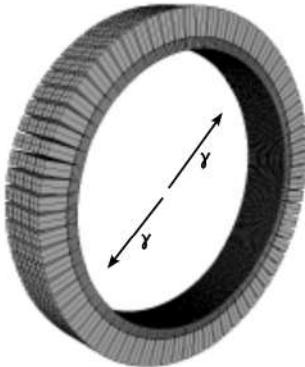
$$\begin{aligned} m &= \frac{E}{c^2} \\ &= \frac{0.5 \times 10^6 \text{ J}}{(3 \times 10^8 \text{ m/s})^2} \\ &= 6 \times 10^{-12} \text{ kilograms.} \end{aligned}$$

The change in mass is too small to measure with any practical technique. This is because the square of the speed of light is such a large number.

Electron-positron annihilation

example 21

Natural radioactivity in the earth produces positrons, which are like electrons but have the opposite charge. A form of antimatter, positrons annihilate with electrons to produce gamma rays, a form of high-frequency light. Such a process would have been considered impossible before Einstein, because conservation of mass and energy were believed to be separate principles, and this process eliminates 100% of the original mass. The amount of energy produced by annihilating 1 kg of matter with 1 kg of



g / Top: A PET scanner. Middle: Each positron annihilates with an electron, producing two gamma-rays that fly off back-to-back. When two gamma rays are observed simultaneously in the ring of detectors, they are assumed to come from the same annihilation event, and the point at which they were emitted must lie on the line connecting the two detectors. Bottom: A scan of a person's torso. The body has concentrated the radioactive tracer around the stomach, indicating an abnormal medical condition.

antimatter is

$$\begin{aligned} E &= mc^2 \\ &= (2 \text{ kg}) (3.0 \times 10^8 \text{ m/s})^2 \\ &= 2 \times 10^{17} \text{ J}, \end{aligned}$$

which is on the same order of magnitude as a day's energy consumption for the entire world's population!

Positron annihilation forms the basis for the medical imaging technique called a PET (positron emission tomography) scan, in which a positron-emitting chemical is injected into the patient and mapped by the emission of gamma rays from the parts of the body where it accumulates.

One commonly hears some misinterpretations of $E = mc^2$, one being that the equation tells us how much kinetic energy an object would have if it was moving at the speed of light. This wouldn't make much sense, both because the equation for kinetic energy has $1/2$ in it, $KE = (1/2)mv^2$, and because a material object can't be made to move at the speed of light. However, this naturally leads to the question of just how much mass-energy a moving object has. We know that when the object is at rest, it has no kinetic energy, so its mass-energy is simply equal to the energy-equivalent of its mass, mc^2 ,

$$\mathcal{E} = mc^2 \text{ when } v = 0,$$

where the symbol \mathcal{E} (cursive "E") stands for mass-energy. The point of using the new symbol is simply to remind ourselves that we're talking about relativity, so an object at rest has $\mathcal{E} = mc^2$, not $E = 0$ as we'd assume in nonrelativistic physics.

Suppose we start accelerating the object with a constant force. A constant force means a constant rate of transfer of momentum, but $p = m\gamma v$ approaches infinity as v approaches c , so the object will only get closer and closer to the speed of light, but never reach it. Now what about the work being done by the force? The force keeps doing work and doing work, which means that we keep on using up energy. Mass-energy is conserved, so the energy being expended must equal the increase in the object's mass-energy. We can continue this process for as long as we like, and the amount of mass-energy will increase without limit. We therefore conclude that an object's mass-energy approaches infinity as its speed approaches the speed of light,

$$\mathcal{E} \rightarrow \infty \text{ when } v \rightarrow c.$$

Now that we have some idea what to expect, what is the actual equation for the mass-energy? As proved in section 7.3.4, it is

$$\mathcal{E} = m\gamma c^2.$$

self-check D

Verify that this equation has the two properties we wanted. ▷

Answer, p. 1062

■ *KE compared to mc^2 at low speeds* *example 22*

▷ An object is moving at ordinary nonrelativistic speeds. Compare its kinetic energy to the energy mc^2 it has purely because of its mass.

▷ The speed of light is a very big number, so mc^2 is a huge number of joules. The object has a gigantic amount of energy because of its mass, and only a relatively small amount of additional kinetic energy because of its motion.

Another way of seeing this is that at low speeds, γ is only a tiny bit greater than 1, so \mathcal{E} is only a tiny bit greater than mc^2 .

■ *The correspondence principle for mass-energy* *example 23*

▷ Show that the equation $\mathcal{E} = m\gamma c^2$ obeys the correspondence principle.

▷ As we accelerate an object from rest, its mass-energy becomes greater than its resting value. Nonrelativistically, we interpret this excess mass-energy as the object's kinetic energy,

$$\begin{aligned} KE &= \mathcal{E}(v) - \mathcal{E}(v = 0) \\ &= m\gamma c^2 - mc^2 \\ &= m(\gamma - 1)c^2. \end{aligned}$$

Expressing γ as $(1 - v^2/c^2)^{-1/2}$ and making use of the approximation $(1 + \epsilon)^p \approx 1 + p\epsilon$ for small ϵ , we have $\gamma \approx 1 + v^2/2c^2$, so

$$\begin{aligned} KE &\approx m(1 + \frac{v^2}{2c^2} - 1)c^2 \\ &= \frac{1}{2}mv^2, \end{aligned}$$

which is the nonrelativistic expression. As demanded by the correspondence principle, relativity agrees with newtonian physics at speeds that are small compared to the speed of light.

7.3.3 * The energy-momentum four-vector

Starting from $\mathcal{E} = m\gamma$ and $p = m\gamma v$, a little algebra allows one to prove the identity

$$m^2 = \mathcal{E}^2 - p^2.$$

We can define an energy-momentum four-vector,

$$\mathbf{p} = (\mathcal{E}, p_x, p_y, p_z),$$

and the relation $m^2 = \mathcal{E}^2 - p^2$ then arises from the inner product $\mathbf{p} \cdot \mathbf{p}$. Since \mathcal{E} and \mathbf{p} are separately conserved, the energy-momentum four-vector is also conserved.

v	Y
0.9870	1.0002(5)
0.9881	1.0012(5)
0.9900	0.9998(5)

A high-precision test of this fundamental relativistic relationship was carried out by Meyer *et al.* in 1963 by studying the motion of electrons in static electric and magnetic fields. They define the quantity

$$Y^2 = \frac{\mathcal{E}^2}{m^2 + p^2},$$

which according to special relativity should equal 1. Their results, tabulated in the sidebar, show excellent agreement with theory.

Energy and momentum of light example 24

Light has $m = 0$ and $\gamma = \infty$, so if we try to apply $\mathcal{E} = m\gamma$ and $p = m\gamma v$ to light, or to any massless particle, we get the indeterminate form $0 \cdot \infty$, which can't be evaluated without a delicate and laborious evaluation of limits as in problem 11 on p. 460.

Applying $m^2 = \mathcal{E}^2 - p^2$ yields the same result, $\mathcal{E} = |p|$, much more easily. This example demonstrates that although we encountered the relations $\mathcal{E} = m\gamma$ and $p = m\gamma v$ first, the identity $m^2 = \mathcal{E}^2 - p^2$ is actually more fundamental.

Figure q on p. 734 shows an experiment that verified $\mathcal{E} = |p|$ empirically.

For the reasons given in example 24, we take $m^2 = \mathcal{E}^2 - p^2$ to be the *definition* of mass in relativity. One thing to be careful about is that this definition is not additive. Suppose that we lump two systems together and call them one big system, adding their mass-energies and momenta. When we do this, the mass of the combination is not the same as the sum of the masses. For example, suppose we have two rays of light moving in opposite directions, with energy-momentum vectors $(\mathcal{E}, \mathcal{E}, 0, 0)$ and $(\mathcal{E}, -\mathcal{E}, 0, 0)$. Adding these gives $(2\mathcal{E}, 0, 0)$, which implies a mass equal to $2\mathcal{E}$. In fact, in the early universe, where the density of light was high, the universe's ambient gravitational fields were mainly those caused by the light it contained.

Mass-energy, not energy, goes in the energy-momentum four-vector example 25

When we say that something is a four-vector, we mean that it behaves properly under a Lorentz transformation: we can draw such a four-vector on graph paper, and then when we change frames of reference, we should be able to measure the vector in the new frame of reference by using the new version of the graph-paper grid derived from the old one by a Lorentz transformation.

If we had used the energy E rather than the mass-energy \mathcal{E} to construct the energy-momentum four-vector, we wouldn't have gotten a valid four-vector. An easy way to see this is to consider the case where a noninteracting object is at rest in some frame of reference. Its momentum and kinetic energy are both zero. If we'd defined $\mathbf{p} = (E, p_x, p_y, p_z)$ rather than $\mathbf{p} = (\mathcal{E}, p_x, p_y, p_z)$, we

would have had $\mathbf{p} = 0$ in this frame. But when we draw a zero vector, we get a point, and a point remains a point regardless of how we distort the graph paper we use to measure it. That wouldn't have made sense, because in other frames of reference, we have $E \neq 0$.

Metric units

example 26

The relation $m^2 = \varepsilon^2 - p^2$ is only valid in relativistic units. If we tried to apply it without modification to numbers expressed in metric units, we would have

$$\text{kg}^2 = \text{kg}^2 \cdot \frac{\text{m}^4}{\text{s}^4} - \text{kg}^2 \cdot \frac{\text{m}^2}{\text{s}^2},$$

which would be nonsense because the three terms all have different units. As usual, we need to insert factors of c to make a metric version, and these factors of c are determined by the need to fix the broken units:

$$m^2 c^4 = \varepsilon^2 - p^2 c^2$$

Pair production requires matter

example 27

Example 21 on p. 435 discussed the annihilation of an electron and a positron into two gamma rays, which is an example of turning matter into pure energy. An opposite example is pair production, a process in which a gamma ray disappears, and its energy goes into creating an electron and a positron.

Pair production cannot happen in a vacuum. For example, gamma rays from distant black holes can travel through empty space for thousands of years before being detected on earth, and they don't turn into electron-positron pairs before they can get here. Pair production can only happen in the presence of matter. When lead is used as shielding against gamma rays, one of the ways the gamma rays can be stopped in the lead is by undergoing pair production.

To see why pair production is forbidden in a vacuum, consider the process in the frame of reference in which the electron-positron pair has zero total momentum. In this frame, the gamma ray would have to have had zero momentum, but a gamma ray with zero momentum must have zero energy as well (example 24). This means that conservation of four-momentum has been violated: the timelike component of the four-momentum is the mass-energy, and it has increased from 0 in the initial state to at least $2mc^2$ in the final state.

7.3.4 ★ Proofs

This optional section proves some results claimed earlier.

Ultrarelativistic motion

We start by considering the case of a particle, described as “ultrarelativistic,” that travels at very close to the speed of light. A good way of thinking about such a particle is that it’s one with a very small mass. For example, the subatomic particle called the neutrino has a very small mass, thousands of times smaller than that of the electron. Neutrinos are emitted in radioactive decay, and because the neutrino’s mass is so small, the amount of energy available in these decays is always enough to accelerate it to very close to the speed of light. Nobody has ever succeeded in observing a neutrino that was *not* ultrarelativistic. When a particle’s mass is very small, the mass becomes difficult to measure. For almost 70 years after the neutrino was discovered, its mass was thought to be zero. Similarly, we currently believe that a ray of light has no mass, but it is always possible that its mass will be found to be nonzero at some point in the future. A ray of light can be modeled as an ultrarelativistic particle.

Let’s compare ultrarelativistic particles with train cars. A single car with kinetic energy E has different properties than a train of two cars each with kinetic energy $E/2$. The single car has half the mass and a speed that is greater by a factor of $\sqrt{2}$. But the same is not true for ultrarelativistic particles. Since an idealized ultrarelativistic particle has a mass too small to be detectable in any experiment, we can’t detect the difference between m and $2m$. Furthermore, ultrarelativistic particles move at close to c , so there is no observable difference in speed. Thus we expect that a single ultrarelativistic particle with energy E compared with two such particles, each with energy $E/2$, should have all the same properties as measured by a mechanical detector.

An idealized zero-mass particle also has no frame in which it can be at rest. It always travels at c , and no matter how fast we chase after it, we can never catch up. We can, however, observe it in different frames of reference, and we will find that its energy is different. For example, distant galaxies are receding from us at substantial fractions of c , and when we observe them through a telescope, they appear very dim not just because they are very far away but also because their light has less energy in our frame than in a frame at rest relative to the source. This effect must be such that changing frames of reference according to a specific Lorentz transformation always changes the energy of the particle by a fixed factor, regardless of the particle’s original energy; for if not, then the effect of a Lorentz transformation on a single particle of energy E would be different from its effect on two particles of energy $E/2$.

How does this energy-shift factor depend on the velocity v of the Lorentz transformation? Rather than v , it becomes more convenient to express things in terms of the Doppler shift factor D , which multiplies when we change frames of reference. Let's write $f(D)$ for the energy-shift factor that results from a given Lorentz transformation. Since a Lorentz transformation D_1 followed by a second transformation D_2 is equivalent to a single transformation by D_1D_2 , we must have $f(D_1D_2) = f(D_1)f(D_2)$. This tightly constrains the form of the function f ; it must be something like $f(D) = s^n$, where n is a constant. We postpone until p. 442 the proof that $n = 1$, which is also in agreement with experiments with rays of light.

Our final result is that the energy of an ultrarelativistic particle is simply proportional to its Doppler shift factor D . Even in the case where the particle is truly massless, so that D doesn't have any finite value, we can still find how the energy differs according to different observers by finding the D of the Lorentz transformation between the two observers' frames of reference.

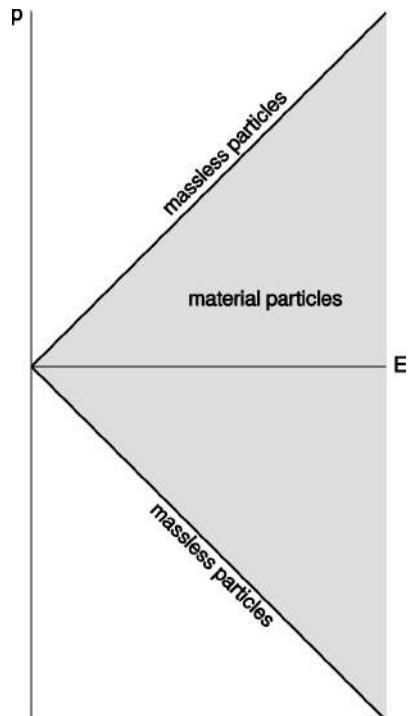
Energy

The following argument is due to Einstein. Suppose that a material object O of mass m , initially at rest in a certain frame A, emits two rays of light, each with energy $E/2$. By conservation of energy, the object must have lost an amount of energy equal to E . By symmetry, O remains at rest.

We now switch to a new frame of reference moving at a certain velocity v in the z direction relative to the original frame. We assume that O 's energy is different in this frame, but that the change in its energy amounts to multiplication by some unitless factor x , which depends only on v , since there is nothing else it could depend on that could allow us to form a unitless quantity. In this frame the light rays have energies $ED(v)$ and $ED(-v)$. If conservation of energy is to hold in the new frame as it did in the old, we must have $2xE = ED(v) + ED(-v)$. After some algebra, we find $x = 1/\sqrt{1 - v^2}$, which we recognize as γ . This proves that $E = m\gamma$ for a material object.

Momentum

We've seen that ultrarelativistic particles are “generic,” in the sense that they have no individual mechanical properties other than an energy and a direction of motion. Therefore the relationship between energy and momentum must be *linear* for ultrarelativistic particles. Indeed, experiments verify that light has momentum, and doubling the energy of a ray of light doubles its momentum rather than quadrupling it. On a graph of p versus E , massless particles, which have $E \propto |p|$, lie on two diagonal lines that connect at the origin. If we like, we can pick units such that the slopes of these lines are plus and minus one. Material particles lie to the right of



In the p - E plane, massless particles lie on the two diagonals, while particles with mass lie to the right.

these lines. For example, a car sitting in a parking lot has $p = 0$ and $E = mc^2$.

Now what happens to such a graph when we change to a different frame or reference that is in motion relative to the original frame? A massless particle still has to act like a massless particle, so the diagonals are simply stretched or contracted along their own lengths. In fact the transformation must be linear (p. 403), because conservation of energy and momentum involve addition, and we need these laws to be valid in all frames of reference. By the same reasoning as in figure j on p. 405, the transformation must be area-preserving. We then have the same three cases to consider as in figure g on p. 404. Case I is ruled out because it would imply that particles keep the same energy when we change frames. (This is what would happen if c were infinite, so that the mass-equivalent E/c^2 of a given energy was zero, and therefore E would be interpreted purely as the mass.) Case II can't be right because it doesn't preserve the $E = |p|$ diagonals. We are left with case III, which establishes the fact that the p - E plane transforms according to exactly the same kind of Lorentz transformation as the x - t plane. That is, (E, p_x, p_y, p_z) is a four-vector.

The only remaining issue to settle is whether the choice of units that gives invariant 45-degree diagonals in the x - t plane is the same as the choice of units that gives such diagonals in the p - E plane. That is, we need to establish that the c that applies to x and t is equal to the c' needed for p and E , i.e., that the velocity scales of the two graphs are matched up. This is true because in the Newtonian limit, the total mass-energy E is essentially just the particle's mass, and then $p/E \approx p/m \approx v$. This establishes that the velocity scales are matched at small velocities, which implies that they coincide for all velocities, since a large velocity, even one approaching c , can be built up from many small increments. (This also establishes that the exponent n defined on p. 441 equals 1 as claimed.)

Since $m^2 = E^2 - p^2$, it follows that for a material particle, $p = m\gamma v$.

7.4 ★ General relativity

What you've learned so far about relativity is known as the special theory of relativity, which is compatible with three of the four known forces of nature: electromagnetism, the strong nuclear force, and the weak nuclear force. Gravity, however, can't be shoehorned into the special theory. In order to make gravity work, Einstein had to generalize relativity. The resulting theory is known as the general theory of relativity.⁷

7.4.1 Our universe isn't Euclidean

Euclid proved thousands of years ago that the angles in a triangle add up to 180° . But what does it really mean to "prove" this? Euclid proved it based on certain assumptions (his five postulates), listed in the margin of this page. But how do we know that the postulates are true?

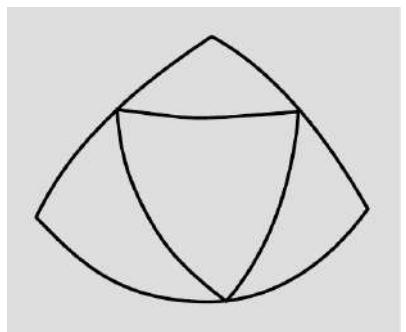
Only by observation can we tell whether any of Euclid's statements are correct characterizations of how space actually behaves in our universe. If we draw a triangle on paper with a ruler and measure its angles with a protractor, we will quickly verify to pretty good precision that the sum is close to 180° . But of course we already knew that space was at least *approximately* Euclidean. If there had been any gross error in Euclidean geometry, it would have been detected in Euclid's own lifetime. The correspondence principle tells us that if there is going to be any deviation from Euclidean geometry, it must be small under ordinary conditions.

To improve the precision of the experiment, we need to make sure that our ruler is very straight. One way to check would be to sight along it by eye, which amounts to comparing its straightness to that of a ray of light. For that matter, we might as well throw the physical ruler in the trash and construct our triangle out of three laser beams. To avoid effects from the air we should do the experiment in outer space. Doing it in space also has the advantage of allowing us to make the triangle very large; as shown in figure a, the discrepancy from 180° is expected to be proportional to the area of the triangle.

But we already know that light rays are bent by gravity. We expect it based on $E = mc^2$, which tells us that the energy of a light ray is equivalent to a certain amount of mass, and furthermore it has been verified experimentally by the deflection of starlight by the sun (example 18, p. 434). We therefore know that our universe is noneuclidean, and we gain the further insight that the level of

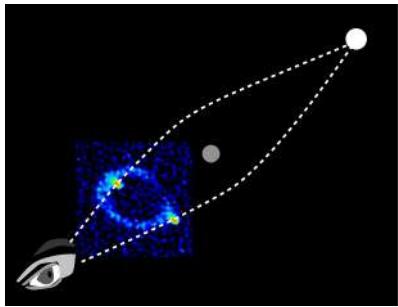
Postulates of Euclidean geometry:

1. Two points determine a line.
2. Line segments can be extended.
3. A unique circle can be constructed given any point as its center and any line segment as its radius.
4. All right angles are equal to one another.
5. Given a line and a point not on the line, no more than one line can be drawn through the point and parallel to the given line.



a / Noneuclidean effects, such as the discrepancy from 180° in the sum of the angles of a triangle, are expected to be proportional to area. Here, a noneuclidean equilateral triangle is cut up into four smaller equilateral triangles, each with $1/4$ the area. As proved in problem 22, the discrepancy is quadrupled when the area is quadrupled.

⁷Einstein originally described the distinction between the two theories by saying that the special theory applied to nonaccelerating frames of reference, while the general one allowed any frame at all. The modern consensus is that Einstein was misinterpreting his own theory, and that special relativity actually handles accelerating frames just fine.



b / An Einstein's ring. The distant object is a quasar, MG1131+0456, and the one in the middle is an unknown object, possibly a supermassive black hole. The intermediate object's gravity focuses the rays of light from the distant one. Because the entire arrangement lacks perfect axial symmetry, the ring is nonuniform; most of its brightness is concentrated in two lumps on opposite sides.

deviation from Euclidean behavior depends on gravity.

Since the noneuclidean effects are bigger when the system being studied is larger, we expect them to be especially important in the study of cosmology, where the distance scales are very large.

Einstein's ring

example 28

An Einstein's ring, figure b, is formed when there is a chance alignment of a distant source with a closer gravitating body. This type of gravitational lensing is direct evidence for the noneuclidean nature of space. The two light rays are lines, and they violate Euclid's first postulate, that two points determine a line.

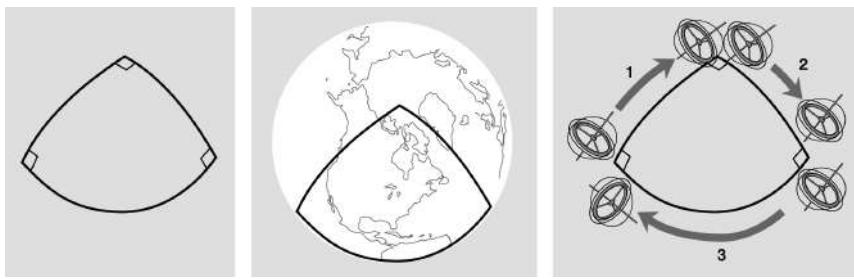
One could protest that effects like these are just an imperfection of the light rays as physical models of straight lines. Maybe the noneuclidean effects would go away if we used something better and straighter than a light ray. But we don't know of anything straighter than a light ray. Furthermore, we observe that all measuring devices, not just optical ones, report the same noneuclidean behavior.

Curvature

An example of such a non-optical measurement is the Gravity Probe B satellite, figure d, which was launched into a polar orbit in 2004 and operated until 2010. The probe carried four gyroscopes made of quartz, which were the most perfect spheres ever manufactured, varying from sphericity by no more than about 40 atoms. Each gyroscope floated weightlessly in a vacuum, so that its rotation was perfectly steady. After 5000 orbits, the gyroscopes had reoriented themselves by about 2×10^{-3} ° relative to the distant stars. This effect cannot be explained by Newtonian physics, since no torques acted on them. It was, however, exactly as predicted by Einstein's theory of general relativity. It becomes easier to see why such an effect would be expected due to the noneuclidean nature of space if we characterize Euclidean geometry as the geometry of a flat plane as opposed to a curved one. On a curved surface like a sphere, figure c, Euclid's fifth postulate fails, and it's not hard to see that we can get triangles for which the sum of the angles is not 180°. By transporting a gyroscope all the way around the edges of such a triangle and back to its starting point, we change its orientation.

The triangle in figure c has angles that add up to more than 180°. This type of curvature is referred to as positive. It is also possible to have negative curvature, as in figure e.

In general relativity, curvature isn't just something caused by gravity. Gravity *is* curvature, and the curvature involves both space and time, as may become clearer once you get to figure k. Thus the distinction between special and general relativity is that general relativity handles curved spacetime, while special relativity is restricted

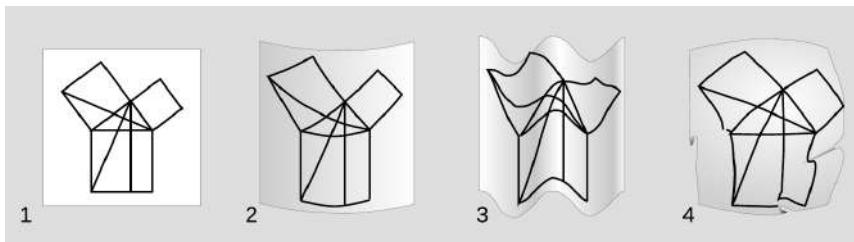


c / Left: A 90-90-90 triangle. Its angles add up to more than 180° . Middle: The triangle “pops” off the page visually. We intuitively want to visualize it as lying on a curved surface such as the earth’s. Right: A gyroscope carried smoothly around its perimeter ends up having changed its orientation when it gets back to its starting point.

to the case where spacetime is flat.

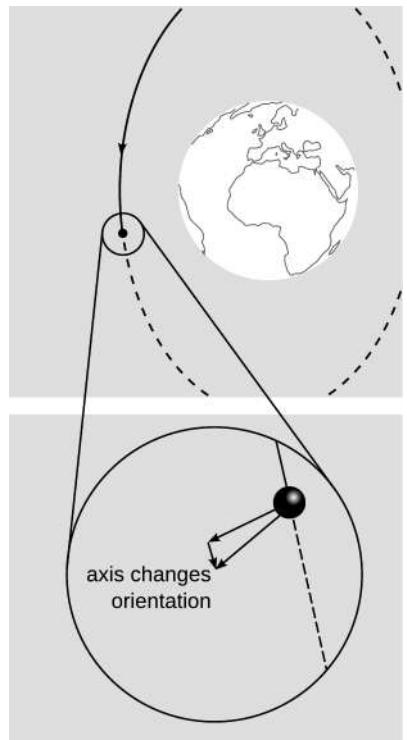
Curvature doesn’t require higher dimensions

Although we often visualize curvature by imagining embedding a two-dimensional surface in a three-dimensional space, that’s just an aid in visualization. There is no evidence for any additional dimensions, nor is it necessary to hypothesize them in order to let spacetime be curved as described in general relativity.

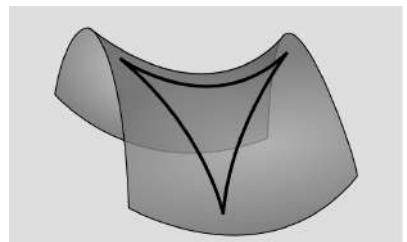


f / Only measurements from within the plane define whether the plane is curved. It could look curved when drawn embedded in three dimensions, but nevertheless still be intrinsically flat.

Put yourself in the shoes of a two-dimensional being living in a two-dimensional space. Euclid’s postulates all refer to constructions that can be performed using a compass and an unmarked straight-edge. If this being can physically verify them all as descriptions of the space she inhabits, then she knows that her space is Euclidean, and that propositions such as the Pythagorean theorem are physically valid in her universe. But the diagram in f/1 illustrating the proof of the Pythagorean theorem in Euclid’s *Elements* (proposition I.47) is equally valid if the page is rolled onto a cylinder, 2, or formed into a wavy corrugated shape, 3. These types of curvature, which can be achieved without tearing or crumpling the surface, are not real to her. They are simply side-effects of visualizing her



d / Gravity Probe B was in a polar orbit around the earth. As in the right panel of figure c, the orientation of the gyroscope changes when it is carried around a curve and back to its starting point. Because the effect was small, it was necessary to let it accumulate over the course of 5000 orbits in order to make it detectable.



e / A triangle in a space with negative curvature has angles that add to less than 180° .

two-dimensional universe as if it were embedded in a hypothetical third dimension — which doesn't exist in any sense that is empirically verifiable to her. Of the curved surfaces in figure f, only the sphere, 4, has curvature that she can measure; the diagram can't be plastered onto the sphere without folding or cutting and pasting.

So the observation of curvature doesn't imply the existence of extra dimensions, nor does embedding a space in a higher-dimensional one so that it looks curvy always mean that there will be any curvature detectable from within the lower-dimensional space.

7.4.2 The equivalence principle

Universality of free-fall

Although light rays and gyroscopes seem to agree that space is curved in a gravitational field, it's always conceivable that we could find something else that would disagree. For example, suppose that there is a new and improved ray called the StraightRayTM. The StraightRay is like a light ray, but when we construct a triangle out of StraightRays, we always get the Euclidean result for the sum of the angles. We would then have to throw away general relativity's whole idea of describing gravity in terms of curvature. One good way of making a StraightRay would be if we had a supply of some kind of exotic matter — call it FloatyStuffTM — that had the ordinary amount of inertia, but was completely unaffected by gravity. We could then shoot a stream of FloatyStuff particles out of a nozzle at nearly the speed of light and make a StraightRay.

Normally when we release a material object in a gravitational field, it experiences a force mg , and then by Newton's second law its acceleration is $a = F/m = mg/m = g$. The m 's cancel, which is the reason that everything falls with the same acceleration (in the absence of other forces such as air resistance). The universality of this behavior is what allows us to interpret the gravity geometrically in general relativity. For example, the Gravity Probe B gyroscopes were made out of quartz, but if they had been made out of something else, it wouldn't have mattered. But if we had access to some FloatyStuff, the geometrical picture of gravity would fail, because the “ m ” that described its susceptibility to gravity would be a different “ m ” than the one describing its inertia.

The question of the existence or nonexistence of such forms of matter turns out to be related to the question of what kinds of motion are relative. Let's say that alien gangsters land in a flying saucer, kidnap you out of your back yard, konk you on the head, and take you away. When you regain consciousness, you're locked up in a sealed cabin in their spaceship. You pull your keychain out of your pocket and release it, and you observe that it accelerates toward the floor with an acceleration that seems quite a bit slower than what you're used to on earth, perhaps a third of a gee. There

are two possible explanations for this. One is that the aliens have taken you to some other planet, maybe Mars, where the strength of gravity is a third of what we have on earth. The other is that your keychain didn't really accelerate at all: you're still inside the flying saucer, which is accelerating at a third of a gee, so that it was really the deck that accelerated up and hit the keys.

There is absolutely no way to tell which of these two scenarios is actually the case — unless you happen to have a chunk of *FloatyStuff* in your other pocket. If you release the *FloatyStuff* and it hovers above the deck, then you're on another planet and experiencing genuine gravity; your keychain responded to the gravity, but the *FloatyStuff* didn't. But if you release the *FloatyStuff* and see it hit the deck, then the flying saucer is accelerating through outer space.

The nonexistence of *FloatyStuff* in our universe is called the *equivalence principle*. If the equivalence principle holds, then an acceleration (such as the acceleration of the flying saucer) is always equivalent to a gravitational field, and no observation can ever tell the difference without reference to something external. (And suppose you did have some external reference point — how would you know whether *it* was accelerating?)

The artificial horizon

example 29

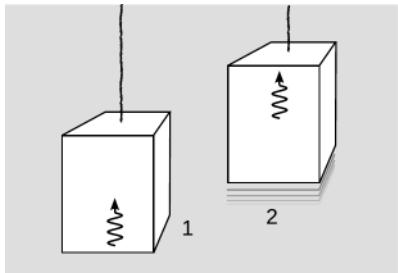
The pilot of an airplane cannot always easily tell which way is up. The horizon may not be level simply because the ground has an actual slope, and in any case the horizon may not be visible if the weather is foggy. One might imagine that the problem could be solved simply by hanging a pendulum and observing which way it pointed, but by the equivalence principle the pendulum cannot tell the difference between a gravitational field and an acceleration of the aircraft relative to the ground — nor can any other accelerometer, such as the pilot's inner ear. For example, when the plane is turning to the right, accelerometers will be tricked into believing that "down" is down and to the left. To get around this problem, airplanes use a device called an artificial horizon, which is essentially a gyroscope. The gyroscope has to be initialized when the plane is known to be oriented in a horizontal plane. No gyroscope is perfect, so over time it will drift. For this reason the instrument also contains an accelerometer, and the gyroscope is always forced into agreement with the accelerometer's average output over the preceding several minutes. If the plane is flown in circles for several minutes, the artificial horizon will be fooled into indicating that the wrong direction is vertical.

Gravitational Doppler shifts and time dilation

An interesting application of the equivalence principle is the explanation of gravitational time dilation. As described on p. 400, experiments show that a clock at the top of a mountain runs faster



g / An artificial horizon.



h / 1. A ray of light is emitted upward from the floor of the elevator. The elevator accelerates upward. 2. By the time the light is detected at the ceiling, the elevator has changed its velocity, so the light is detected with a Doppler shift.



i / Pound and Rebka at the top and bottom of the tower.

than one down at its foot.

To calculate this effect, we make use of the fact that the gravitational field in the area around the mountain is equivalent to an acceleration. Suppose we're in an elevator accelerating upward with acceleration a , and we shoot a ray of light from the floor up toward the ceiling, at height h . The time Δt it takes the light ray to get to the ceiling is about h/c , and by the time the light ray reaches the ceiling, the elevator has sped up by $v = a\Delta t = ah/c$, so we'll see a red-shift in the ray's frequency. Since v is small compared to c , we don't need to use the fancy Doppler shift equation from subsection 7.2.8; we can just approximate the Doppler shift factor as $1 - v/c \approx 1 - ah/c^2$. By the equivalence principle, we should expect that if a ray of light starts out low down and then rises up through a gravitational field g , its frequency will be Doppler shifted by a factor of $1 - gh/c^2$. This effect was observed in a famous experiment carried out by Pound and Rebka in 1959. Gamma-rays were emitted at the bottom of a 22.5-meter tower at Harvard and detected at the top with the Doppler shift predicted by general relativity. (See problem 25.)

In the mountain-valley experiment, the frequency of the clock in the valley therefore appears to be running too slowly by a factor of $1 - gh/c^2$ when it is compared via radio with the clock at the top of the mountain. We conclude that time runs more slowly when one is lower down in a gravitational field, and the slow-down factor between two points is given by $1 - gh/c^2$, where h is the difference in height.

We have built up a picture of light rays interacting with gravity. To confirm that this make sense, recall that we have already observed in subsection 7.3.3 and in problem 11 on p. 460 that light has momentum. The equivalence principle says that whatever has inertia must also participate in gravitational interactions. Therefore light waves must have weight, and must lose energy when they rise through a gravitational field.

Local flatness

The noneuclidean nature of spacetime produces effects that grow in proportion to the area of the region being considered. Interpreting such effects as evidence of curvature, we see that this connects naturally to the idea that curvature is undetectable from close up. For example, the curvature of the earth's surface is not normally noticeable to us in everyday life. Locally, the earth's surface is flat, and the same is true for spacetime.

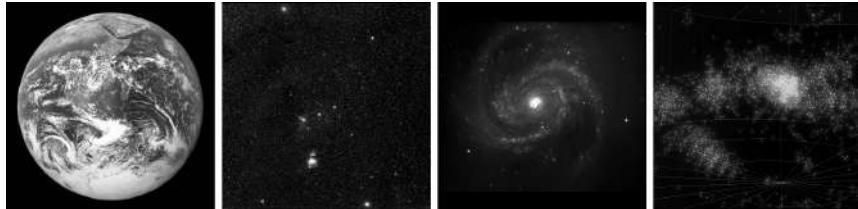
Local flatness turns out to be another way of stating the equivalence principle. In a variation on the alien-abduction story, suppose that you regain consciousness aboard the flying saucer and find yourself weightless. If the equivalence principle holds, then

you have no way of determining from local observations, inside the saucer, whether you are actually weightless in deep space, or simply free-falling in apparent weightlessness, like the astronauts aboard the International Space Station. That means that locally, we can always adopt a free-falling frame of reference in which there is no gravitational field at all. If there is no gravity, then special relativity is valid, and we can treat our local region of spacetime as being approximately flat.

In figure k, an apple falls out of a tree. Its path is a “straight” line in spacetime, in the same sense that the equator is a “straight” line on the earth’s surface.

Inertial frames

In Newtonian mechanics, we have a distinction between inertial and noninertial frames of reference. An inertial frame according to Newton is one that has a constant velocity vector relative to the stars. But what if the stars themselves are accelerating due to a gravitational force from the rest of the galaxy? We could then take the galaxy’s center of mass as defining an inertial frame, but what if something else is acting on the galaxy?



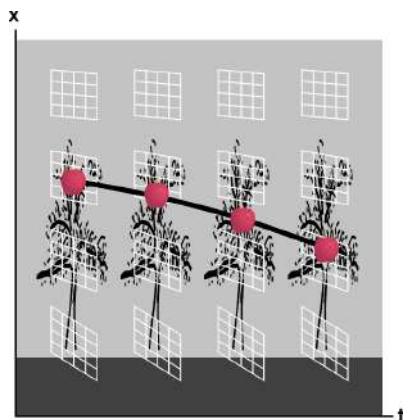
I / Wouldn’t it be nice if we could define the meaning of a Newtonian inertial frame of reference? Newton makes it sound easy: to define an inertial frame, just find some object that is not accelerating because it is not being acted on by any external forces. But what object would we use? The earth? The “fixed stars?” Our galaxy? Our supercluster of galaxies? All of these are accelerating — relative to something.

If we had some *FloatyStuff*, we could resolve the whole question. *FloatyStuff* isn’t affected by gravity, so if we release a sample of it in mid-air, it will continue on a trajectory that defines a perfect Newtonian inertial frame. (We’d better have it on a tether, because otherwise the earth’s rotation will carry the earth out from under it.) But if the equivalence principle holds, then Newton’s definition of an inertial frame is fundamentally flawed.

There is a different definition of an inertial frame that works better in relativity. A Newtonian inertial frame was defined by an object that isn’t subject to any forces, gravitational or otherwise. In general relativity, we instead define an inertial frame using an



j / The earth is flat — locally.



k / Spacetime is locally flat.

object that that isn't influenced by anything other than gravity. By this definition, a free-falling rock defines an inertial frame, but this book sitting on your desk does not.

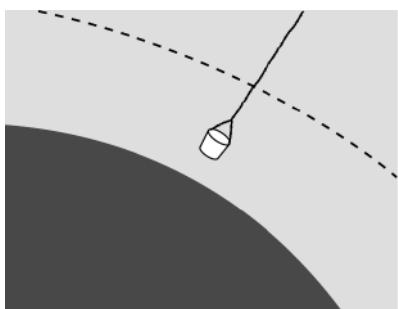
7.4.3 Black holes

The observations described so far showed only small effects from curvature. To get a big effect, we should look at regions of space in which there are strong gravitational fields. The prime example is a black hole. The best studied examples are two objects in our own galaxy: Cygnus X-1, which is believed to be a black hole with about ten times the mass of our sun, and Sagittarius A*, an object near the center of our galaxy with about four million solar masses.

Although a black hole is a relativistic object, we can gain some insight into how it works by applying Newtonian physics. A spherical body of mass M has an escape velocity $v = \sqrt{2GM/r}$, which is the minimum velocity that we would need to give to a projectile shot from a distance r so that it would never fall back down. If r is small enough, the escape velocity will be greater than c , so that even a ray of light can never escape.

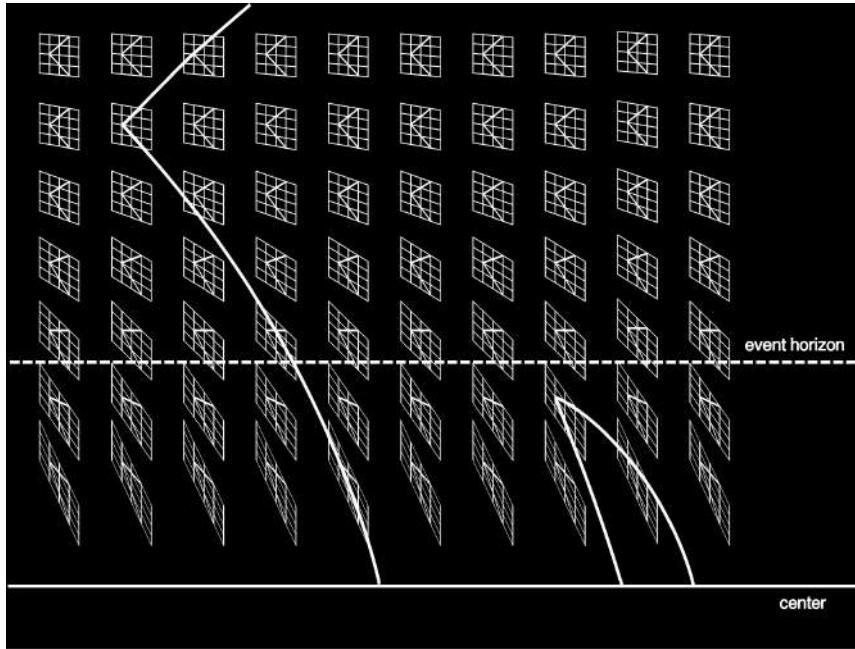
We can now make an educated guess as to what this means without having to study all the mathematics of general relativity. In relativity, c isn't really the speed of light, it's really to be thought of as a restriction on how fast cause and effect can propagate through space. This suggests the correct interpretation, which is that for an object compact enough to be a black hole, there is no way for an event at a distance closer than r to have an effect on an event far away. There is an invisible, spherical boundary with radius r , called the event horizon, and the region within that boundary is cut off from the rest of the universe in terms of cause and effect. If you wanted to explore that region, you could drop into it while wearing a space-suit — but it would be a one-way trip, because you could never get back out to report on what you had seen.

In the Newtonian description of a black hole, matter could be lifted out of a black hole, m . Would this be possible with a real-world black hole, which is relativistic rather than Newtonian? No, because the bucket is causally separated from the outside universe. No rope would be strong enough for this job (problem 12, p. 460).



m / Matter is lifted out of a Newtonian black hole with a bucket. The dashed line represents the point at which the escape velocity equals the speed of light. For a real, relativistic black hole, this is impossible.

One misleading aspect of the Newtonian analysis is that it encourages us to imagine that a light ray trying to escape from a black hole will slow down, stop, and then fall back in. This can't be right, because we know that any observer who sees a light ray flying by always measures its speed to be c . This was true in special relativity, and by the equivalence principle we can be assured that the same is true *locally* in general relativity. Figure n shows what would really happen.



Although the light rays in figure n don't speed up or slow down, they do experience gravitational Doppler shifts. If a light ray is emitted from just above the event horizon, then it will escape to an infinite distance, but it will suffer an extreme Doppler shift toward low frequencies. A distant observer also has the option of interpreting this as a gravitational time dilation that greatly lowers the frequency of the oscillating electric charges that produced the ray. If the point of emission is made closer and closer to the horizon, the frequency and energy measured by a distant observer approach zero, making the ray impossible to observe.

Information paradox

Black holes have some disturbing implications for the kind of universe that in the Age of the Enlightenment was imagined to have been set in motion initially and then left to run forever like clockwork.

Newton's laws have built into them the implicit assumption that omniscience is possible, at least in principle. For example, Newton's definition of an inertial frame of reference leads to an infinite regress, as described on p. 449. For Newton this isn't a problem, because in principle an omniscient observer can know the location of every mass in the universe. In this conception of the cosmos, there are no theoretical limits on human knowledge, only practical ones; if we could gather sufficiently precise data about the state of the universe at one time, and if we could carry out all the calculations to extrapolate into the future, then we could know everything that would ever happen. (See the famous quote by Laplace on p. 16.)

But the existence of event horizons surrounding black holes makes

n / The equivalence principle tells us that spacetime locally has the same structure as in special relativity, so we can draw the familiar parallelogram of $x - t$ coordinates at each point near the black hole. Superimposed on each little grid is a pair of lines representing motion at the speed of light in both directions, inward and outward. Because spacetime is curved, these lines do not appear to be at 45-degree angles, but to an observer in that region, they would appear to be. When light rays are emitted inward and outward from a point outside the event horizon, one escapes and one plunges into the black hole. On this diagram, they look like they are decelerating and accelerating, but local observers comparing them to their own coordinate grids would always see them as moving at exactly c . When rays are emitted from a point *inside* the event horizon, neither escapes; the distortion is so severe that "outward" is really inward.

it impossible for any observer to be omniscient; only an observer inside a particular horizon can see what's going on inside that horizon.

Furthermore, a black hole has at its center an infinitely dense point, called a singularity, containing all its mass, and this implies that information can be destroyed and made inaccessible to *any* observer at all. For example, suppose that astronaut Alice goes on a suicide mission to explore a black hole, free-falling in through the event horizon. She has a certain amount of time to collect data and satisfy her intellectual curiosity, but then she impacts the singularity and is compacted into a mathematical point. Now astronaut Betty decides that she will never be satisfied unless the secrets revealed to Alice are known to her as well — and besides, she was Alice's best friend, and she wants to know whether Alice had any last words. Betty can jump through the horizon, but she can never know Alice's last words, nor can any other observer who jumps in after Alice does.

This destruction of information is known as the black hole information paradox, and it's referred to as a paradox because quantum physics (ch. 13) has built into its DNA the requirement that information is never lost in this sense.

Formation

Around 1960, as black holes and their strange properties began to be better understood and more widely discussed, many physicists who found these issues distressing comforted themselves with the belief that black holes would never really form from realistic initial conditions, such as the collapse of a massive star. Their skepticism was not entirely unreasonable, since it is usually very hard in astronomy to hit a gravitating target, the reason being that conservation of angular momentum tends to make the projectile swing past. (See problem 13 on p. 295 for a quantitative analysis.) For example, if we wanted to drop a space probe into the sun, we would have to extremely precisely stop its sideways orbital motion so that it would drop almost exactly straight in. Once a star started to collapse, the theory went, and became relatively compact, it would be such a small target that further infalling material would be unlikely to hit it, and the process of collapse would halt. According to this point of view, theorists who had calculated the collapse of a star into a black hole had been oversimplifying by assuming a star that was initially perfectly spherical and nonrotating. Remove the unrealistically perfect symmetry of the initial conditions, and a black hole would never actually form.

But Roger Penrose proved in 1964 that this was wrong. In fact, once an object collapses to a certain density, the Penrose singularity theorem guarantees mathematically that it will collapse further until a singularity is formed, and this singularity is surrounded by an event horizon. Since the brightness of an object like Sagittarius A* is far too low to be explained unless it has an event horizon (the



o / In Newtonian contexts, physicists and astronomers had a correct intuition that it's hard for things to collapse gravitationally. This star cluster has been around for about 15 billion years, but it hasn't collapsed into a black hole. If any individual star happens to head toward the center, conservation of angular momentum tends to cause it to swing past and fly back out. The Penrose singularity theorem tells us that this Newtonian intuition is wrong when applied to an object that has collapsed past a certain point.

interstellar gas flowing into it would glow due to frictional heating), we can be certain that there really is a singularity at its core.

7.4.4 Cosmology

The Big Bang

Subsection 6.1.5 presented the evidence, discovered by Hubble, that the universe is expanding in the aftermath of the Big Bang: when we observe the light from distant galaxies, it is always Doppler-shifted toward the red end of the spectrum, indicating that no matter what direction we look in the sky, everything is rushing away from us. This seems to go against the modern attitude, originated by Copernicus, that we and our planet do not occupy a special place in the universe. Why is everything rushing away from *our* planet in particular? But general relativity shows that this anti-Copernican conclusion is wrong. General relativity describes space not as a rigidly defined background but as something that can curve and stretch, like a sheet of rubber. We imagine all the galaxies as existing on the surface of such a sheet, which then expands uniformly. The space between the galaxies (but not the galaxies themselves) grows at a steady rate, so that any observer, inhabiting any galaxy, will see every other galaxy as receding. There is therefore no privileged or special location in the universe.

We might think that there would be another kind of special place, which would be the one at which the Big Bang happened. Maybe someone has put a brass plaque there? But general relativity doesn't describe the Big Bang as an explosion that suddenly occurred in a preexisting background of time and space. According to general relativity, space itself came into existence at the Big Bang, and the hot, dense matter of the early universe was uniformly distributed everywhere. The Big Bang happened everywhere at once.

Observations show that the universe is very uniform on large scales, and for ease of calculation, the first physical models of the expanding universe were constructed with perfect uniformity. In these models, the Big Bang was a singularity. This singularity can't even be included as an event in spacetime, so that time itself only exists after the Big Bang. A Big Bang singularity also creates an even more acute version of the black hole information paradox. Whereas matter and information disappear *into* a black hole singularity, stuff pops *out* of a Big Bang singularity, and there is no physical principle that could predict what it would be.

As with black holes, there was considerable skepticism about whether the existence of an initial singularity in these models was an artifact of the unrealistically perfect uniformity assumed in the models. Perhaps in the real universe, extrapolation of all the paths of the galaxies backward in time would show them missing each other by millions of light-years. But in 1972 Stephen Hawking proved

a variant on the Penrose singularity theorem that applied to Big Bang singularities. By the Hawking singularity theorem, the level of uniformity we see in the present-day universe is more than sufficient to prove that a Big Bang singularity must have existed.

The cosmic censorship hypothesis

It might not be too much of a philosophical jolt to imagine that information was spontaneously created in the Big Bang. Setting up the initial conditions of the entire universe is traditionally the prerogative of God, not the laws of physics. But there is nothing fundamental in general relativity that forbids the existence of other singularities that act like the Big Bang, being information producers rather than information consumers. As John Earman of the University of Pittsburgh puts it, anything could pop out of such a singularity, including green slime or your lost socks. This would eliminate any hope of finding a universal set of laws of physics that would be able to make a prediction given any initial situation.

That would be such a devastating defeat for the enterprise of physics that in 1969 Penrose proposed an alternative, humorously named the “cosmic censorship hypothesis,” which states that every singularity in our universe, other than the Big Bang, is hidden behind an event horizon. Therefore if green slime spontaneously pops out of one, there is limited impact on the predictive ability of physics, since the slime can never have any causal effect on the outside world. A singularity that is not modestly cloaked behind an event horizon is referred to as a naked singularity. Nobody has yet been able to prove the cosmic censorship hypothesis.

The advent of high-precision cosmology

We expect that if there is matter in the universe, it should have gravitational fields, and in the rubber-sheet analogy this should be represented as a curvature of the sheet. Instead of a flat sheet, we can have a spherical balloon, so that cosmological expansion is like inflating it with more and more air. It is also possible to have negative curvature, as in figure e on p. 445. All three of these are valid, possible cosmologies according to relativity. The positive-curvature type happens if the average density of matter in the universe is above a certain critical level, the negative-curvature one if the density is below that value.

To find out which type of universe we inhabit, we could try to take a survey of the matter in the universe and determine its average density. Historically, it has been very difficult to do this, even to within an order of magnitude. Most of the matter in the universe probably doesn’t emit light, making it difficult to detect. Astronomical distance scales are also very poorly calibrated against absolute units such as the SI.

Instead, we measure the universe’s curvature, and infer the den-



p / An expanding universe with positive spatial curvature can be imagined as a balloon being blown up. Every galaxy's distance from every other galaxy increases, but no galaxy is the center of the expansion.

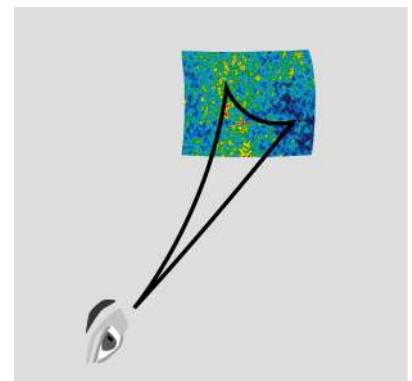
sity of matter from that. It turns out that we can do this by observing the cosmic microwave background (CMB) radiation, which is the light left over from the brightly glowing early universe, which was dense and hot. As the universe has expanded, light waves that were in flight have expanded their wavelengths along with it. This afterglow of the big bang was originally visible light, but after billions of years of expansion it has shifted into the microwave radio part of the electromagnetic spectrum. The CMB is not perfectly uniform, and this turns out to give us a way to measure the universe's curvature. Since the CMB was emitted when the universe was only about 400,000 years old, any vibrations or disturbances in the hot hydrogen and helium gas that filled space in that era would only have had time to travel a certain distance, limited by the speed of sound. We therefore expect that no feature in the CMB should be bigger than a certain known size. In a universe with negative spatial curvature, the sum of the interior angles of a triangle is less than the Euclidean value of 180 degrees. Therefore if we observe a variation in the CMB over some angle, the distance between two points on the sky is actually greater than would have been inferred from Euclidean geometry. The opposite happens if the curvature is positive.

This observation was done by the 1989-1993 COBE probe, and its 2001-2009 successor, the Wilkinson Microwave Anisotropy Probe. The result is that the angular sizes are almost exactly *equal* to what they should be according to Euclidean geometry. We therefore infer that the universe is very close to having zero average spatial curvature on the cosmological scale, and this tells us that its average density must be within about 0.5% of the critical value. The years since COBE and WMAP mark the advent of an era in which cosmology has gone from being a field of estimates and rough guesses to a high-precision science.

If one is inclined to be skeptical about the seemingly precise answers to the mysteries of the cosmos, there are consistency checks that can be carried out. In the bad old days of low-precision cosmology, estimates of the age of the universe ranged from 10 billion to 20 billion years, and the low end was inconsistent with the age of the oldest star clusters. This was believed to be a problem either for observational cosmology or for the astrophysical models used to estimate the ages of the clusters: “You can’t be older than your ma.” Current data have shown that the low estimates of the age were incorrect, so consistency is restored. (The best figure for the age of the universe is currently 13.8 ± 0.1 billion years.)

Dark energy and dark matter

Not everything works out so smoothly, however. One surprise is that the universe's expansion is not currently slowing down, as had been expected due to the gravitational attraction of all the matter



q / The angular scale of fluctuations in the cosmic microwave background can be used to infer the curvature of the universe.

in it. Instead, it is currently speeding up. This is attributed to a variable in Einstein’s equations, long assumed to be zero, which represents a universal gravitational repulsion of space itself, occurring even when there is no matter present. The current name for this is “dark energy,” although the fancy name is just a label for our ignorance about what causes it.

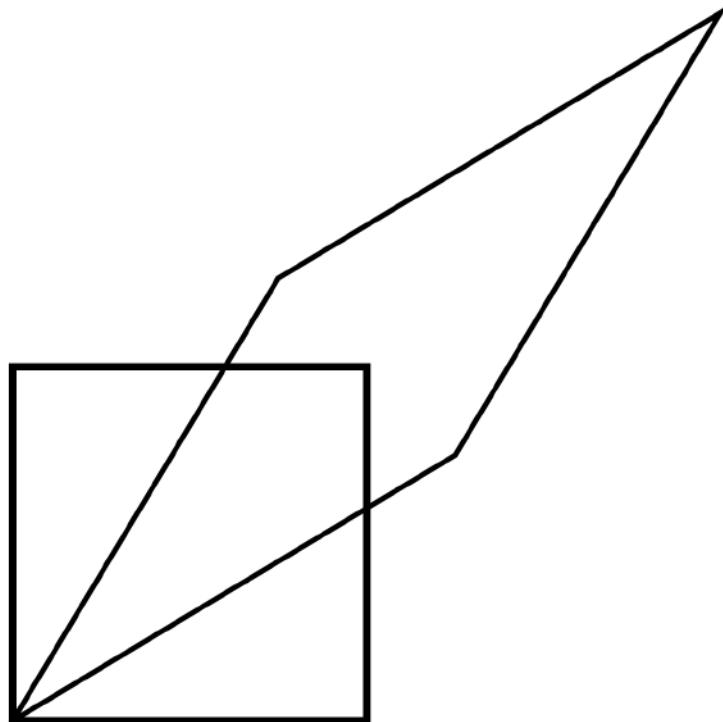
Another surprise comes from attempts to model the formation of the elements during the era shortly after the Big Bang, before the formation of the first stars. The observed relative abundances of hydrogen, helium, and deuterium (^2H) cannot be reconciled with the density of low-velocity matter inferred from the observational data. If the inferred mass density were entirely due to normal matter (i.e., matter whose mass consisted mostly of protons and neutrons), then nuclear reactions in the dense early universe should have proceeded relatively efficiently, leading to a much higher ratio of helium to hydrogen, and a much lower abundance of deuterium. The conclusion is that most of the matter in the universe must be made of an unknown type of exotic matter, known as “dark matter.” We are in the ironic position of knowing that precisely 96% of the universe is something other than atoms, but knowing nothing about what that something is. As of 2013, there have been several experiments that have been carried out to attempt the direct detection of dark matter particles. These are carried out at the bottom of mineshafts to eliminate background radiation. Early claims of success appear to have been statistical flukes, and the most sensitive experiments have not detected anything.⁸

⁸arxiv.org/abs/1310.8214

Problems

The symbols \checkmark , \blacksquare , etc. are explained on page 464.

- 1 The figure illustrates a Lorentz transformation using the conventions employed in section 7.2. For simplicity, the transformation chosen is one that lengthens one diagonal by a factor of 2. Since Lorentz transformations preserve area, the other diagonal is shortened by a factor of 2. Let the original frame of reference, depicted with the square, be A, and the new one B. (a) By measuring with a ruler on the figure, show that the velocity of frame B relative to frame A is $0.6c$. (b) Print out a copy of the page. With a ruler, draw a third parallelogram that represents a second successive Lorentz transformation, one that lengthens the long diagonal by another factor of 2. Call this third frame C. Use measurements with a ruler to determine frame C's velocity relative to frame A. Does it equal double the velocity found in part a? Explain why it should be expected to turn out the way it does. \checkmark \blacksquare



2 Astronauts in three different spaceships are communicating with each other. Those aboard ships A and B agree on the rate at which time is passing, but they disagree with the ones on ship C.
(a) Alice is aboard ship A. How does she describe the motion of her own ship, in its frame of reference?

- (b) Describe the motion of the other two ships according to Alice.
(c) Give the description according to Betty, whose frame of reference is ship B.
(d) Do the same for Cathy, aboard ship C.

3 What happens in the equation for γ when you put in a negative number for v ? Explain what this means physically, and why it makes sense.

4 The Voyager 1 space probe, launched in 1977, is moving faster relative to the earth than any other human-made object, at 17,000 meters per second.

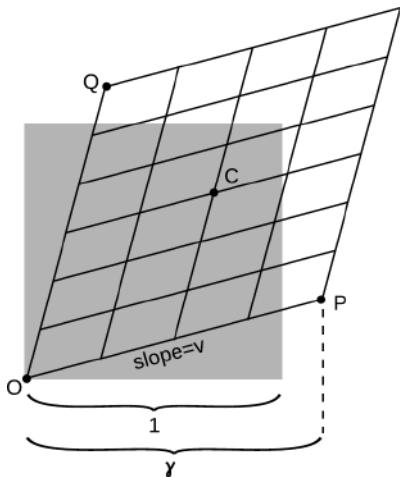
- (a) Calculate the probe's γ .
(b) Over the course of one year on earth, slightly less than one year passes on the probe. How much less? (There are 31 million seconds in a year.)

5 In example 5 on page 408, I remarked that accelerating a macroscopic (i.e., not microscopic) object to close to the speed of light would require an unreasonable amount of energy. Suppose that the starship Enterprise from Star Trek has a mass of 8.0×10^7 kg, about the same as the Queen Elizabeth 2. Compute the kinetic energy it would have if it was moving at half the speed of light. Compare with the total energy content of the world's nuclear arsenals, which is about 10^{21} J.

✓

6 The earth is orbiting the sun, and therefore is contracted relativistically in the direction of its motion. Compute the amount by which its diameter shrinks in this direction.

✓



Problem 7.

7 In this homework problem, you'll fill in the steps of the algebra required in order to find the equation for γ on page 405. To keep the algebra simple, let the time t in figure k equal 1, as suggested in the figure accompanying this homework problem. The original square then has an area of 1, and the transformed parallelogram must also have an area of 1.
(a) Prove that point P is at $x = v\gamma$, so that its (t, x) coordinates are $(\gamma, v\gamma)$.
(b) Find the (t, x) coordinates of point Q.
(c) Find the length of the short diagonal connecting P and Q.
(d) Average the coordinates of P and Q to find the coordinates of the midpoint C of the parallelogram, and then find distance OC.
(e) Find the area of the parallelogram by computing twice the area of triangle PQO. [Hint: You can take PQ to be the base of the triangle.]
(f) Set this area equal to 1 and solve for γ to prove $\gamma = 1/\sqrt{1 - v^2}$.

✓

- 8** (a) A free neutron (as opposed to a neutron bound into an atomic nucleus) is unstable, and undergoes beta decay (which you may want to review). The masses of the particles involved are as follows:

neutron	1.67495×10^{-27} kg
proton	1.67265×10^{-27} kg
electron	0.00091×10^{-27} kg
antineutrino	$< 10^{-35}$ kg

Find the energy released in the decay of a free neutron. ✓

(b) Neutrons and protons make up essentially all of the mass of the ordinary matter around us. We observe that the universe around us has no free neutrons, but lots of free protons (the nuclei of hydrogen, which is the element that 90% of the universe is made of). We find neutrons only inside nuclei along with other neutrons and protons, not on their own.

If there are processes that can convert neutrons into protons, we might imagine that there could also be proton-to-neutron conversions, and indeed such a process does occur sometimes in nuclei that contain both neutrons and protons: a proton can decay into a neutron, a positron, and a neutrino. A positron is a particle with the same properties as an electron, except that its electrical charge is positive. A neutrino, like an antineutrino, has negligible mass.

Although such a process can occur within a nucleus, explain why it cannot happen to a free proton. (If it could, hydrogen would be radioactive, and you wouldn't exist!) ■

- 9** (a) Find a relativistic equation for the velocity of an object in terms of its mass and momentum (eliminating γ). Use natural units (i.e., discard factors of c) throughout. ✓

(b) Show that your result is approximately the same as the nonrelativistic value, p/m , at low velocities.

(c) Show that very large momenta result in speeds close to the speed of light.

(d) Insert factors of c to make your result from part a usable in SI units. ✓



- 10** (a) Show that for $v = (3/5)c$, γ comes out to be a simple fraction.

(b) Find another value of v for which γ is a simple fraction. ■

11 An object moving at a speed very close to the speed of light is referred to as ultrarelativistic. Ordinarily (luckily) the only ultrarelativistic objects in our universe are subatomic particles, such as cosmic rays or particles that have been accelerated in a particle accelerator.

- (a) What kind of number is γ for an ultrarelativistic particle?
- (b) Repeat example 22 on page 437, but instead of very low, non-relativistic speeds, consider ultrarelativistic speeds.
- (c) Find an equation for the ratio \mathcal{E}/p . The speed may be relativistic, but don't assume that it's ultrarelativistic. ✓
- (d) Simplify your answer to part c for the case where the speed is ultrarelativistic. ✓
- (e) We can think of a beam of light as an ultrarelativistic object — it certainly moves at a speed that's sufficiently close to the speed of light! Suppose you turn on a one-watt flashlight, leave it on for one second, and then turn it off. Compute the momentum of the recoiling flashlight, in units of $\text{kg}\cdot\text{m/s}$. ✓
- (f) Discuss how part e relates to the correspondence principle.

12 As discussed in chapter 6, the speed at which a disturbance travels along a string under tension is given by $v = \sqrt{T/\mu}$, where μ is the mass per unit length, and T is the tension.

- (a) Suppose a string has a density ρ , and a cross-sectional area A . Find an expression for the maximum tension that could possibly exist in the string without producing $v > c$, which is impossible according to relativity. Express your answer in terms of ρ , A , and c . The interpretation is that relativity puts a limit on how strong any material can be. ✓
- (b) Every substance has a tensile strength, defined as the force per unit area required to break it by pulling it apart. The tensile strength is measured in units of N/m^2 , which is the same as the pascal (Pa), the mks unit of pressure. Make a numerical estimate of the maximum tensile strength allowed by relativity in the case where the rope is made out of ordinary matter, with a density on the same order of magnitude as that of water. (For comparison, kevlar has a tensile strength of about $4 \times 10^9 \text{ Pa}$, and there is speculation that fibers made from carbon nanotubes could have values as high as $6 \times 10^{10} \text{ Pa}$.) ✓
- (c) A black hole is a star that has collapsed and become very dense, so that its gravity is too strong for anything ever to escape from it. For instance, the escape velocity from a black hole is greater than c , so a projectile can't be shot out of it. Many people, when they hear this description of a black hole in terms of an escape velocity, wonder why it still wouldn't be possible to extract an object from a black hole by other means. For example, suppose we lower an astronaut into a black hole on a rope, and then pull him back out again. Why might this not work?

- 13** (a) A charged particle is surrounded by a uniform electric field. Starting from rest, it is accelerated by the field to speed v after traveling a distance d . Now it is allowed to continue for a further distance $3d$, for a total displacement from the start of $4d$. What speed will it reach, assuming newtonian physics?
 (b) Find the relativistic result for the case of $v = c/2$.

14 Problem 14 has been deleted.

- 15** Expand the equation $K = m(\gamma - 1)$ in a Taylor series, and find the first two nonvanishing terms. Explain why the vanishing terms are the ones that should vanish physically. Show that the first term is the nonrelativistic expression for kinetic energy.

- 16** Consider the relativistic relation for momentum as a function of velocity (for a particle with nonzero mass). Expand this in a Taylor series, and find the first two nonvanishing terms. Explain why the vanishing terms are the ones that should vanish physically. Show that the first term is the newtonian expression.

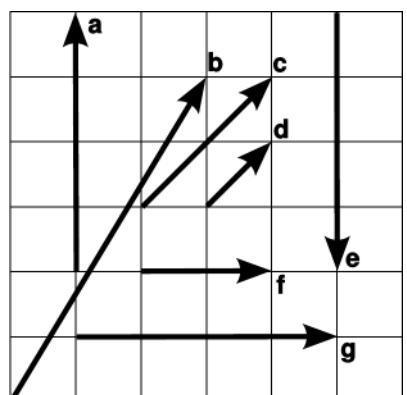
- 17** As promised in subsection 7.2.8, this problem will lead you through the steps of finding an equation for the combination of velocities in relativity, generalizing the numerical result found in problem 1. Suppose that A moves relative to B at velocity u , and B relative to C at v . We want to find A's velocity w relative to C, in terms of u and v . Suppose that A emits light with a certain frequency. This will be observed by B with a Doppler shift $D(u)$. C detects a further shift of $D(v)$ relative to B. We therefore expect the Doppler shifts to multiply, $D(w) = D(u)D(v)$, and this provides an implicit rule for determining w if u and v are known. (a) Using the expression for D given in section 7.2.8, write down an equation relating u , v , and w . (b) Solve for w in terms of u and v . (c) Show that your answer to part b satisfies the correspondence principle.

▷ Solution, p. 1046

- 18** The figure shows seven four-vectors, represented in a two-dimensional plot of x versus t . All the vectors have y and z components that are zero. Which of these vectors are congruent to others, i.e., which represent spacetime intervals that are equal to one another? If you reason based on Euclidean geometry, you will get the wrong answers.

▷ Solution, p. 1047

- 19** Four-vectors can be timelike, lightlike, or spacelike. What can you say about the inherent properties of particles whose momentum four-vectors fall in these various categories?



Problem 18.

20 The following are the three most common ways in which gamma rays interact with matter:

Photoelectric effect: The gamma ray hits an electron, is annihilated, and gives all of its energy to the electron.

Compton scattering: The gamma ray bounces off of an electron, exiting in some direction with some amount of energy.

Pair production: The gamma ray is annihilated, creating an electron and a positron.

Example 27 on p. 439 shows that pair production can't occur in a vacuum due to conservation of the energy-momentum four-vector. What about the other two processes? Can the photoelectric effect occur without the presence of some third particle such as an atomic nucleus? Can Compton scattering happen without a third particle?

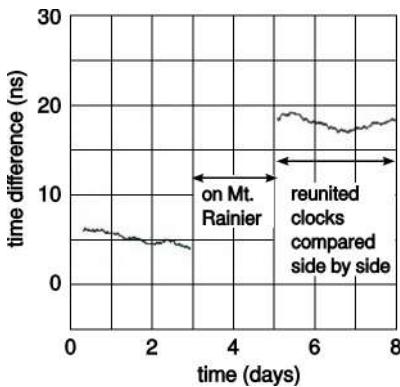
21 Expand the relativistic equation for the longitudinal Doppler shift of light $D(v)$ in a Taylor series, and find the first two nonvanishing terms. Show that these two terms agree with the nonrelativistic expression, so that any relativistic effect is of higher order in v . ■

22 Prove, as claimed in the caption of figure a on p. 443, that $S - 180^\circ = 4(s - 180^\circ)$, where S is the sum of the angles of the large equilateral triangle and s is the corresponding sum for one of the four small ones. ▷ Solution, p. 1047 ■

23 If a two-dimensional being lived on the surface of a cone, would it say that its space was curved, or not? ■

24 (a) Verify that the equation $1 - gh/c^2$ for the gravitational Doppler shift and gravitational time dilation has units that make sense. (b) Does this equation satisfy the correspondence principle?

25 (a) Calculate the Doppler shift to be expected in the Pound-Rebka experiment described on p. 448. (b) In the 1978 Iijima mountain-valley experiment (p. 400), analysis was complicated by the clock's sensitivity to pressure, humidity, and temperature. A cleaner version of the experiment was done in 2005 by hobbyist Tom Van Baak. He put his kids and three of his atomic clocks in a minivan and drove from Bellevue, Washington to a lodge on Mount Rainier, 1340 meters higher in elevation. They spent the weekend there. Back at home, he compared the clocks to others that had stayed at his house. Verify that the effect shown in the graph is as predicted by general relativity. ■



Problem 25b. Redrawn from Van Baak, Physics Today 60 (2007) 16.

26 The International Space Station orbits at an altitude of about 350 km and a speed of about 8000 m/s relative to the ground. Compare the gravitational and kinematic time dilations. Over all, does time run faster on the ISS than on the ground, or more slowly? ■

27 Section 7.4.3 presented a Newtonian estimate of how compact an object would have to be in order to be a black hole. Although this estimate is not really right, it turns out to give the right answer to within about a factor of 2. To roughly what size would the earth have to be compressed in order to become a black hole? ■

28 Clock A sits on a desk. Clock B is tossed up in the air from the same height as the desk and then comes back down. Compare the elapsed times. ▷ Hint, p. 1036 ▷ Solution, p. 1047 ■

29 The angular defect d of a triangle (measured in radians) is defined as $s - \pi$, where s is the sum of the interior angles. The angular defect is proportional to the area A of the triangle. Consider the geometry measured by a two-dimensional being who lives on the surface of a sphere of radius R . First find some triangle on the sphere whose area and angular defect are easy to calculate. Then determine the general equation for d in terms of A and R . ✓ ■

30 (a) In this chapter we've represented Lorentz transformations as distortions of a square into various parallelograms, with the degree of distortion depending on the velocity of one frame of reference relative to another. Suppose that one frame of reference was moving at c relative to another. Discuss what would happen in terms of distortion of a square, and show that this is impossible by using an argument similar to the one used to rule out transformations like the one in figure e on page 403.

(b) Resolve the following paradox. Two pen-pointer lasers are placed side by side and aimed in parallel directions. Their beams both travel at c relative to the hardware, but each beam has a velocity of zero relative to the neighboring beam. But the speed of light can't be zero; it's supposed to be the same in all frames of reference. ■

31 The products of a certain radioactive decay are a massive particle and a gamma ray, which is massless. See example 24 on p. 438 for a discussion of the energy and momentum of a gamma ray. (a) Show that, in the center of mass frame, the energy of the gamma is less than the mass-energy of the massive particle.

(b) Show that the opposite inequality holds if we compare the *kinetic* energy of the massive particle to the energy of the gamma. [Problem by B. Shotwell.] ■

32 Natural relativistic units were introduced on p. 406, and examples 1 and 2 on pp. 407 and 407 gave examples of how to convert an equation from natural units to SI units. In example 4 on p. 408, we derived the approximation

$$\gamma \approx 1 + \frac{v^2}{2}$$

for values of v that are small compared to 1 (i.e., small compared to the speed of light in natural units). As in the other examples, convert this equation to SI units. \triangleright Solution, p. 1047 ■

33 We want to throw a ball of diameter b through a hole of diameter h in a thin wall. Clearly this is possible if $b < h$, but consider the case where $b > h$. If the motion is relativistic, then is it unambiguous whether the ball fits through the hole, or is this frame-dependent, as in example 7 on p. 410? If the former, then is there some velocity v that is required, expressible in terms of b and h ? ■

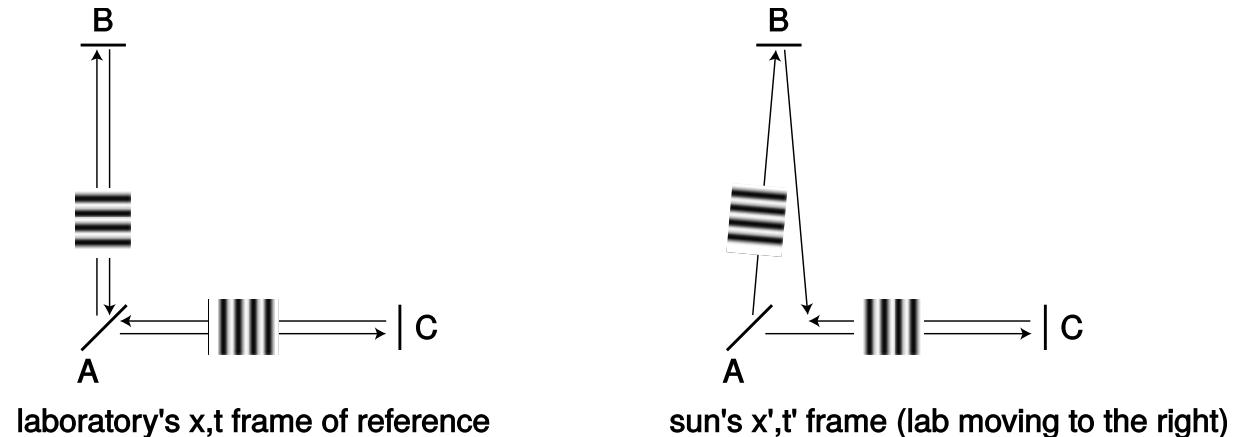
34 (a) Let L be the diameter of our galaxy. Suppose that a person in a spaceship of mass m wants to travel across the galaxy at constant speed, taking proper time τ . Find the kinetic energy of the spaceship. (b) Your friend is impatient, and wants to make the voyage in an hour. For $L = 10^5$ light years, estimate the energy in units of megatons of TNT (1 megaton = 4×10^9 J). ■

Key to symbols:

■ easy ■ typical ■ challenging ■ difficult ■ very difficult
✓ An answer check is available at www.lightandmatter.com.

Exercises

Exercise 7A: The Michelson-Morley Experiment



In this exercise you will analyze the Michelson-Morley experiment, and find what the results should have been according to Galilean relativity and Einstein's theory of relativity. A beam of light coming from the west (not shown) comes to the half-silvered mirror A. Half the light goes through to the east, is reflected by mirror C, and comes back to A. The other half is reflected north by A, is reflected by B, and also comes back to A. When the beams reunite at A, part of each ends up going south, and these parts interfere with one another. If the time taken for a round trip differs by, for example, half the period of the wave, there will be destructive interference.

The point of the experiment was to search for a difference in the experimental results between the daytime, when the laboratory was moving west relative to the sun, and the nighttime, when the laboratory was moving east relative to the sun. Galilean relativity and Einstein's theory of relativity make different predictions about the results. According to Galilean relativity, the speed of light cannot be the same in all reference frames, so it is assumed that there is one special reference frame, perhaps the sun's, in which light travels at the same speed in all directions; in other frames, Galilean relativity predicts that the speed of light will be different in different directions, e.g., slower if the observer is chasing a beam of light. There are four different ways to analyze the experiment:

- *Laboratory's frame of reference, Galilean relativity.* This is not a useful way to analyze the experiment, since one does not know how fast light will travel in various directions.
- *Sun's frame of reference, Galilean relativity.* We assume that in this special frame of reference, the speed of light is the same in all directions: we call this speed c . In this frame, the laboratory moves with velocity v , and mirrors A, B, and C move while the light beam is in flight.
- *Laboratory's frame of reference, Einstein's theory of relativity.* The analysis is extremely simple. Let the length of each arm be L . Then the time required to get from A to either mirror is L/c , so each beam's round-trip time is $2L/c$.
- *Sun's frame of reference, Einstein's theory of relativity.* We analyze this case by starting with the laboratory's frame of reference and then transforming to the sun's frame.

Groups 1-4 work in the sun's frame of reference according to Galilean relativity.

Group 1 finds time AC. Group 2 finds time CA. Group 3 finds time AB. Group 4 finds time BA.

Groups 5 and 6 transform the lab-frame results into the sun's frame according to Einstein's theory.

Group 5 transforms the x and t when ray ACA gets back to A into the sun's frame of reference, and group 6 does the same for ray ABA.

Discussion:

Michelson and Morley found no change in the interference of the waves between day and night. Which version of relativity is consistent with their results?

What does each theory predict if v approaches c ?

What if the arms are not exactly equal in length?

Does it matter if the “special” frame is some frame other than the sun’s?

Exercise 7B: Sports in Slowlightland

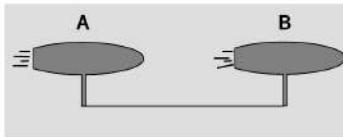
In Slowlightland, the speed of light is $20 \text{ mi/hr} \approx 32 \text{ km/hr} \approx 9 \text{ m/s}$. Think of an example of how relativistic effects would work in sports. Things can get very complex very quickly, so try to think of a simple example that focuses on just one of the following effects:

- relativistic momentum
- relativistic kinetic energy
- relativistic addition of velocities
- time dilation and length contraction
- Doppler shifts of light
- equivalence of mass and energy
- time it takes for light to get to an athlete's eye
- deflection of light rays by gravity

Exercise 7C: Events and Spacetime

Bell's spaceship paradox

A difficult philosophical question is whether the time dilation and length contractions predicted by relativity are "real." This depends, of course, on what one means by "real." They are frame-dependent, i.e., observers in different frames of reference disagree about them. But this doesn't tell us much about their reality, since velocities are frame-dependent in Newtonian mechanics, but nobody worries about whether velocities are real. John Bell (1928-1990) proposed the following thought experiment to physicists in the CERN cafeteria, and found that nearly all of them got it wrong. He took this as evidence that their intuitions had been misguided by the standard way of approaching this question of the reality of Lorentz contractions.



Let spaceships A and B accelerate as shown in the figure, along a straight line. Observer C, whose frame of reference is indicated by the square, does not accelerate. The accelerations, as judged by C, are constant, and equal for the two ships. Each ship is equipped with a yard-arm, and a thread is tied between the two arms. Does the thread break, due to length contraction?

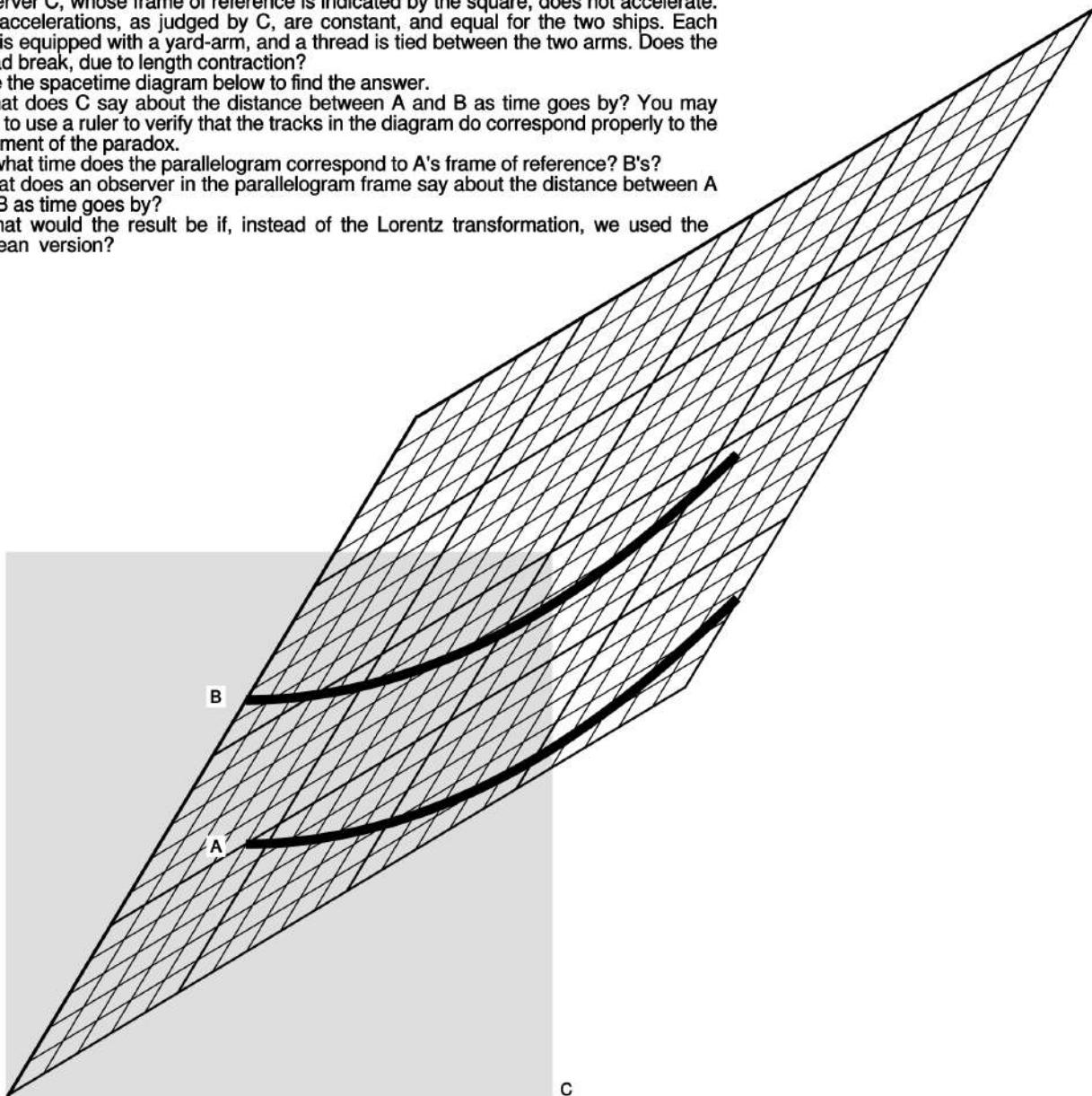
Use the spacetime diagram below to find the answer.

What does C say about the distance between A and B as time goes by? You may want to use a ruler to verify that the tracks in the diagram do correspond properly to the statement of the paradox.

At what time does the parallelogram correspond to A's frame of reference? B's?

What does an observer in the parallelogram frame say about the distance between A and B as time goes by?

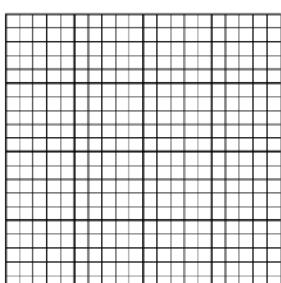
What would the result be if, instead of the Lorentz transformation, we used the Galilean version?



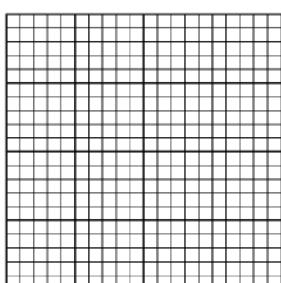
The following paper-and-pencil exercises involve graphs of position versus time. These are often referred to as spacetime diagrams, and tracks on them are called world-lines. As in the other examples in this book, a time axis is always closer to horizontal, and a position axis closer to vertical, and the units are such that the universal speed is a slope of +1 or -1.

Draw spacetime diagrams of: (1) a box with a particle bouncing back and forth inside it, (2) a ray of light being absorbed by a leaf on a tree, (3) an atomic nucleus splitting up into two parts (fissioning), (4) a cloud of gas collapsing gravitationally to form our solar system

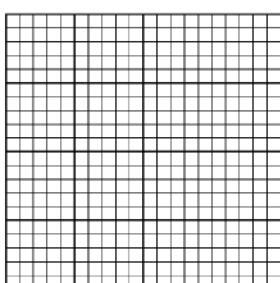
1



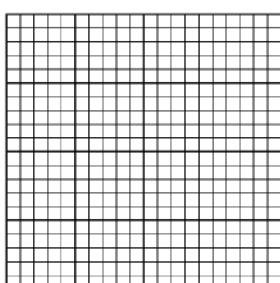
2



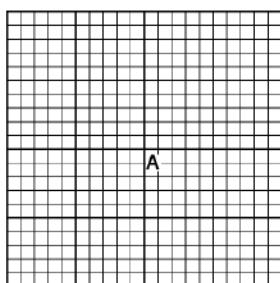
3



4



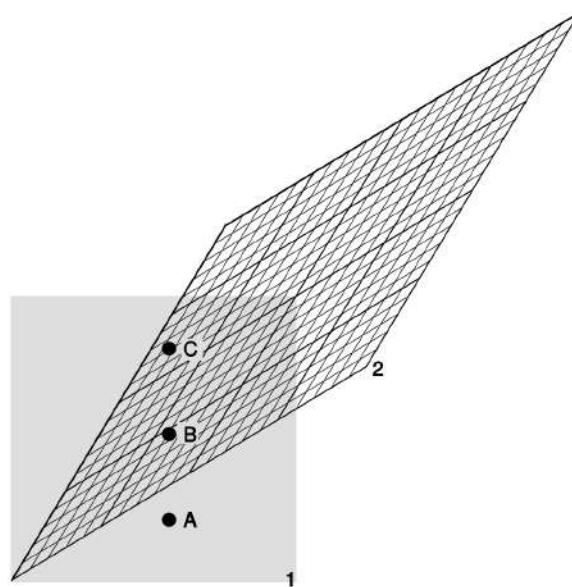
5. Event A is given. Mark examples of events B-F that satisfy these criteria. A material object travels from A to B. An object traveled from C to A. A ray of light emitted at A is received at D. A ray of light emitted at E is received at A. F can have no possible cause-and-effect relationship with A.



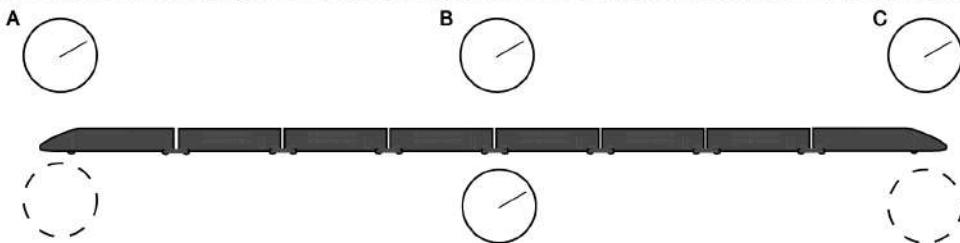
6. These are five pairs of spacetime diagrams, organized so that there is some resemblance between the left and right of each pair. In each case, decide whether the two could actually be the same thing represented in a different inertial frame of reference. Write yes or no, and explain why.



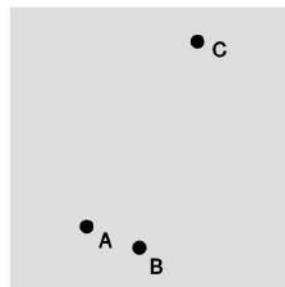
7. The diagram shows three events and two frames of reference. Describe the time-ordering of the events in these frames. Is any other time-ordering possible in any other frame?



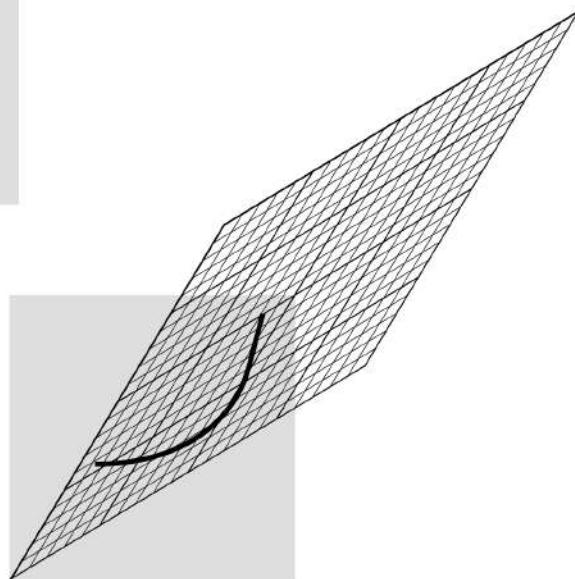
8. The top row shows three clocks located in three different places. They have been synchronized in the frame of reference of the earth, represented by the paper. This synchronization is carried out by exchanging light signals. For example, if the front and back clocks both send out flashes of light when they think it's 2 o'clock, the one in the middle will receive them both at the same time. Event A is the one at which the back clock A reads 2 o'clock, etc. The second row represents clocks that are synchronized aboard the train, which is moving to the right at a substantial fraction of the speed of light. How should the clocks shown with dashed outlines compare with the one at the middle of the train?



9. The figure shows three events and a square representing the t and x axes of a frame of reference. In this frame, the time-ordering of the events is ABC. If we switch to another frame, what other orderings, if any, are possible?



10. The figure shows the motion of an object, with two different frames of reference superimposed. Is anything strange going on?



11. Suppose that the train in 8 is hooked up in a circle, like a chain necklace. What happens?

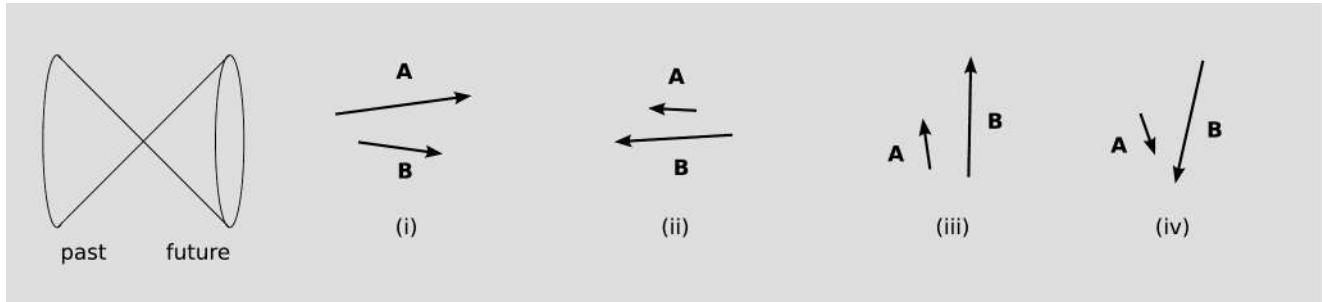
Exercise 7D: Misconceptions about Relativity

The following is a list of common misconceptions about relativity. The class will be split up into random groups, and each group will cooperate on developing an explanation of the misconception, and then the groups will present their explanations to the class. There may be multiple rounds, with students assigned to different randomly chosen groups in successive rounds.

1. How can light have momentum if it has zero mass?
2. What does the world look like in a frame of reference moving at c ?
3. Alice observes Betty coming toward her from the left at $c/2$, and Carol from the right at $c/2$. Therefore Betty is moving at the speed of light relative to Carol.
4. Are relativistic effects such as length contraction and time dilation real, or do they just seem to be that way?
5. Special relativity only matters if you're moving close to the speed of light.
6. Special relativity says that everything is relative.
7. There is a common misconception that relativistic length contraction is what we would actually *see*. Refute this by drawing a spacetime diagram for an object approaching an observer, and tracing rays of light emitted from the object's front and back that both reach the observer's eye at the same time.
8. When you travel close to the speed of light, your time slows down.
9. Is a light wave's wavelength relativistically length contracted by a factor of gamma?
10. Accelerate a baseball to ultrarelativistic speeds. Does it become a black hole?
11. Where did the Big Bang happen?
12. The universe can't be infinite in size, because it's only had a finite amount of time to expand from the point where the Big Bang happened.

Exercise 7E: The sum of observer-vectors is an observer-vector.

The figure gives four pairs of four-vectors, oriented in our customary way as shown by the light-cone on the left.



1. Of the types shown in the four cases i-iv, which types of vectors could represent the world-line of an observer?
2. Suppose that \mathbf{U} and \mathbf{V} are both observer-vectors. What would it mean physically to compute $\mathbf{U} + \mathbf{V}$?
3. Determine the sign of each inner product $\mathbf{A} \cdot \mathbf{B}$.
4. Given an observer whose world-line is along a four-vector \mathbf{O} , suppose we want to determine whether some other four-vector \mathbf{P} is also a possible world-line of an observer. Show that knowledge of the signs of the inner products $\mathbf{O} \cdot \mathbf{P}$ and $\mathbf{P} \cdot \mathbf{P}$ is necessary and sufficient to determine this. Hint: Consider various possibilities like i-iv for vector \mathbf{P} , and see how the signs would turn out.
5. For vectors as described in 4, determine the signs of

$$(\mathbf{U} + \mathbf{V}) \cdot (\mathbf{U} + \mathbf{V})$$

and

$$(\mathbf{U} + \mathbf{V}) \cdot \mathbf{U}$$

by multiplying them out. Interpret the result physically.

Chapter 8

Atoms and Electromagnetism



8.1 The electric glue

Where the telescope ends, the microscope begins. Which of the two has the grander view?

Victor Hugo

His father died during his mother's pregnancy. Rejected by her as a boy, he was packed off to boarding school when she remarried. He himself never married, but in middle age he formed an intense relationship with a much younger man, a relationship that he terminated when he underwent a psychotic break. Following his early scientific successes, he spent the rest of his professional life mostly in frustration over his inability to unlock the secrets of alchemy.

The man being described is Isaac Newton, but not the triumphant Newton of the standard textbook hagiography. Why dwell on the sad side of his life? To the modern science educator, Newton's life-long obsession with alchemy may seem an embarrassment, a distraction from his main achievement, the creation of the modern science of mechanics. To Newton, however, his alchemical researches were naturally related to his investigations of force and motion. What was radical about Newton's analysis of motion was its universality: it succeeded in describing both the heavens and the earth with the same equations, whereas previously it had been assumed that

the sun, moon, stars, and planets were fundamentally different from earthly objects. But Newton realized that if science was to describe all of nature in a unified way, it was not enough to unite the human scale with the scale of the universe: he would not be satisfied until he fit the microscopic universe into the picture as well.

It should not surprise us that Newton failed. Although he was a firm believer in the existence of atoms, there was no more experimental evidence for their existence than there had been when the ancient Greeks first posited them on purely philosophical grounds. Alchemy labored under a tradition of secrecy and mysticism. Newton had already almost single-handedly transformed the fuzzyheaded field of “natural philosophy” into something we would recognize as the modern science of physics, and it would be unjust to criticize him for failing to change alchemy into modern chemistry as well. The time was not ripe. The microscope was a new invention, and it was cutting-edge science when Newton’s contemporary Hooke discovered that living things were made out of cells.

8.1.1 The quest for the atomic force

Newton was not the first of the age of reason. He was the last of the magicians.
John Maynard Keynes

Newton's quest

Nevertheless it will be instructive to pick up Newton's train of thought and see where it leads us with the benefit of modern hindsight. In uniting the human and cosmic scales of existence, he had reimagined both as stages on which the actors were objects (trees and houses, planets and stars) that interacted through attractions and repulsions. He was already convinced that the objects inhabiting the microworld were atoms, so it remained only to determine what kinds of forces they exerted on each other.

His next insight was no less brilliant for his inability to bring it to fruition. He realized that the many human-scale forces — friction, sticky forces, the normal forces that keep objects from occupying the same space, and so on — must all simply be expressions of a more fundamental force acting between atoms. Tape sticks to paper because the atoms in the tape attract the atoms in the paper. My house doesn't fall to the center of the earth because its atoms repel the atoms of the dirt under it.

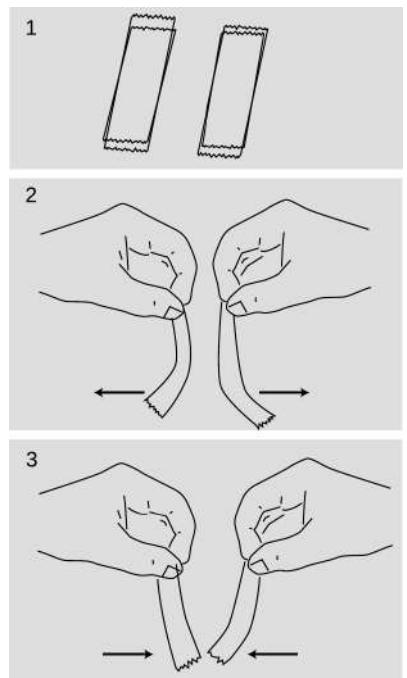
Here he got stuck. It was tempting to think that the atomic force was a form of gravity, which he knew to be universal, fundamental, and mathematically simple. Gravity, however, is always attractive, so how could he use it to explain the existence of both attractive and repulsive atomic forces? The gravitational force between objects of ordinary size is also extremely small, which is why we never notice cars and houses attracting us gravitationally. It would be hard to understand how gravity could be responsible for anything

as vigorous as the beating of a heart or the explosion of gunpowder. Newton went on to write a million words of alchemical notes filled with speculation about some other force, perhaps a “divine force” or “vegetative force” that would for example be carried by the sperm to the egg.

Luckily, we now know enough to investigate a different suspect as a candidate for the atomic force: electricity. Electric forces are often observed between objects that have been prepared by rubbing (or other surface interactions), for instance when clothes rub against each other in the dryer. A useful example is shown in figure a/1: stick two pieces of tape on a tabletop, and then put two more pieces on top of them. Lift each pair from the table, and then separate them. The two top pieces will then repel each other, a/2, as will the two bottom pieces. A bottom piece will attract a top piece, however, a/3. Electrical forces like these are similar in certain ways to gravity, the other force that we already know to be fundamental:

- Electrical forces are *universal*. Although some substances, such as fur, rubber, and plastic, respond more strongly to electrical preparation than others, all matter participates in electrical forces to some degree. There is no such thing as a “nonelectric” substance. Matter is both inherently gravitational and inherently electrical.
- Experiments show that the electrical force, like the gravitational force, is an *inverse square* force. That is, the electrical force between two spheres is proportional to $1/r^2$, where r is the center-to-center distance between them.

Furthermore, electrical forces make more sense than gravity as candidates for the fundamental force between atoms, because we have observed that they can be either attractive or repulsive.



a / Four pieces of tape are prepared, 1, as described in the text. Depending on which combination is tested, the interaction can be either repulsive, 2, or attractive, 3.

8.1.2 Charge, electricity and magnetism

Charge

“Charge” is the technical term used to indicate that an object has been prepared so as to participate in electrical forces. This is to be distinguished from the common usage, in which the term is used indiscriminately for anything electrical. For example, although we speak colloquially of “charging” a battery, you may easily verify that a battery has no charge in the technical sense, e.g., it does not exert any electrical force on a piece of tape that has been prepared as described in the previous section.

Two types of charge

We can easily collect reams of data on electrical forces between different substances that have been charged in different ways. We find for example that cat fur prepared by rubbing against rabbit fur will attract glass that has been rubbed on silk. How can we make any sense of all this information? A vast simplification is achieved by noting that there are really only two types of charge. Suppose we pick cat fur rubbed on rabbit fur as a representative of type A, and glass rubbed on silk for type B. We will now find that there is no “type C.” Any object electrified by any method is either A-like, attracting things A attracts and repelling those it repels, or B-like, displaying the same attractions and repulsions as B. The two types, A and B, always display opposite interactions. If A displays an attraction with some charged object, then B is guaranteed to undergo repulsion with it, and vice-versa.

The coulomb

Although there are only two types of charge, each type can come in different amounts. The metric unit of charge is the coulomb (rhymes with “drool on”), defined as follows:

One Coulomb (C) is the amount of charge such that a force of 9.0×10^9 N occurs between two pointlike objects with charges of 1 C separated by a distance of 1 m.

The notation for an amount of charge is q . The numerical factor in the definition is historical in origin, and is not worth memorizing. The definition is stated for pointlike, i.e., very small, objects, because otherwise different parts of them would be at different distances from each other.

A model of two types of charged particles

Experiments show that all the methods of rubbing or otherwise charging objects involve two objects, and both of them end up getting charged. If one object acquires a certain amount of one type of charge, then the other ends up with an equal amount of the other type. Various interpretations of this are possible, but the simplest is that the basic building blocks of matter come in two flavors, one with each type of charge. Rubbing objects together results in the transfer of some of these particles from one object to the other. In this model, an object that has not been electrically prepared may actually possesses a great deal of *both* types of charge, but the amounts are equal and they are distributed in the same way throughout it. Since type A repels anything that type B attracts, and vice versa, the object will make a total force of zero on any other object. The rest of this chapter fleshes out this model and discusses how these mysterious particles can be understood as being internal parts of atoms.

Use of positive and negative signs for charge

Because the two types of charge tend to cancel out each other's forces, it makes sense to label them using positive and negative signs, and to discuss the *total* charge of an object. It is entirely arbitrary which type of charge to call negative and which to call positive. Benjamin Franklin decided to describe the one we've been calling "A" as negative, but it really doesn't matter as long as everyone is consistent with everyone else. An object with a total charge of zero (equal amounts of both types) is referred to as electrically *neutral*.

self-check A

Criticize the following statement: "There are two types of charge, attractive and repulsive." ▷ Answer, p.

1062

A large body of experimental observations can be summarized as follows:

Coulomb's law: The magnitude of the force acting between point-like charged objects at a center-to-center distance r is given by the equation

$$|\mathbf{F}| = k \frac{|q_1||q_2|}{r^2},$$

where the constant k equals $9.0 \times 10^9 \text{ N}\cdot\text{m}^2/\text{C}^2$. The force is attractive if the charges are of different signs, and repulsive if they have the same sign.

Clever modern techniques have allowed the $1/r^2$ form of Coulomb's law to be tested to incredible accuracy, showing that the exponent is in the range from 1.999999999999998 to 2.0000000000000002.

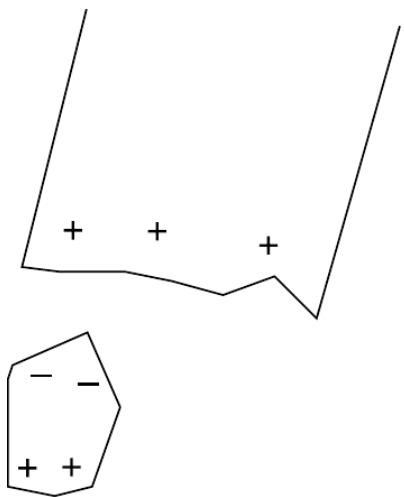
Note that Coulomb's law is closely analogous to Newton's law of gravity, where the magnitude of the force is Gm_1m_2/r^2 , except that there is only one type of mass, not two, and gravitational forces are never repulsive. Because of this close analogy between the two types of forces, we can recycle a great deal of our knowledge of gravitational forces. For instance, there is an electrical equivalent of the shell theorem: the electrical forces exerted externally by a uniformly charged spherical shell are the same as if all the charge was concentrated at its center, and the forces exerted internally are zero.

Conservation of charge

An even more fundamental reason for using positive and negative signs for electrical charge is that experiments show that charge is conserved according to this definition: in any closed system, the total amount of charge is a constant. This is why we observe that rubbing initially uncharged substances together always has the result that one gains a certain amount of one type of charge, while



b / A charged piece of tape attracts uncharged pieces of paper from a distance, and they leap up to it.



c / The paper has zero total charge, but it does have charged particles in it that can move.

the other acquires an equal amount of the other type. Conservation of charge seems natural in our model in which matter is made of positive and negative particles. If the charge on each particle is a fixed property of that type of particle, and if the particles themselves can be neither created nor destroyed, then conservation of charge is inevitable.

Electrical forces involving neutral objects

As shown in figure b, an electrically charged object can attract objects that are uncharged. How is this possible? The key is that even though each piece of paper has a total charge of zero, it has at least some charged particles in it that have some freedom to move. Suppose that the tape is positively charged, c. Mobile particles in the paper will respond to the tape's forces, causing one end of the paper to become negatively charged and the other to become positive. The attraction between the paper and the tape is now stronger than the repulsion, because the negatively charged end is closer to the tape.

self-check B

What would have happened if the tape was negatively charged? ▶

Answer, p. 1062

The path ahead

We have begun to encounter complex electrical behavior that we would never have realized was occurring just from the evidence of our eyes. Unlike the pulleys, blocks, and inclined planes of mechanics, the actors on the stage of electricity and magnetism are invisible phenomena alien to our everyday experience. For this reason, the flavor of the second half of your physics education is dramatically different, focusing much more on experiments and techniques. Even though you will never actually see charge moving through a wire, you can learn to use an ammeter to measure the flow.

Students also tend to get the impression from their first semester of physics that it is a dead science. Not so! We are about to pick up the historical trail that leads directly to the cutting-edge physics research you read about in the newspaper. The atom-smashing experiments that began around 1900, which we will be studying in this chapter, were not that different from the ones of the year 2000 — just smaller, simpler, and much cheaper.

Magnetic forces

A detailed mathematical treatment of magnetism won't come until much later in this book, but we need to develop a few simple ideas about magnetism now because magnetic forces are used in the experiments and techniques we come to next. Everyday magnets

come in two general types. Permanent magnets, such as the ones on your refrigerator, are made of iron or substances like steel that contain iron atoms. (Certain other substances also work, but iron is the cheapest and most common.) The other type of magnet, an example of which is the ones that make your stereo speakers vibrate, consist of coils of wire through which electric charge flows. Both types of magnets are able to attract iron that has not been magnetically prepared, for instance the door of the refrigerator.

A single insight makes these apparently complex phenomena much simpler to understand: magnetic forces are interactions between moving charges, occurring in addition to the electric forces. Suppose a permanent magnet is brought near a magnet of the coiled-wire type. The coiled wire has moving charges in it because we force charge to flow. The permanent magnet also has moving charges in it, but in this case the charges that naturally swirl around inside the iron. (What makes a magnetized piece of iron different from a block of wood is that the motion of the charge in the wood is random rather than organized.) The moving charges in the coiled-wire magnet exert a force on the moving charges in the permanent magnet, and vice-versa.

The mathematics of magnetism is significantly more complex than the Coulomb force law for electricity, which is why we will wait until chapter 11 before delving deeply into it. Two simple facts will suffice for now:

(1) If a charged particle is moving in a region of space near where other charged particles are also moving, their magnetic force on it is directly proportional to its velocity.

(2) The magnetic force on a moving charged particle is always perpendicular to the direction the particle is moving.

A magnetic compass

example 1

The Earth is molten inside, and like a pot of boiling water, it roils and churns. To make a drastic oversimplification, electric charge can get carried along with the churning motion, so the Earth contains moving charge. The needle of a magnetic compass is itself a small permanent magnet. The moving charge inside the earth interacts magnetically with the moving charge inside the compass needle, causing the compass needle to twist around and point north.

A television tube

example 2

A TV picture is painted by a stream of electrons coming from the back of the tube to the front. The beam scans across the whole surface of the tube like a reader scanning a page of a book. Magnetic forces are used to steer the beam. As the beam comes from the back of the tube to the front, up-down and left-right forces are needed for steering. But magnetic forces cannot be used

to get the beam up to speed in the first place, since they can only push perpendicular to the electrons' direction of motion, not forward along it.

Discussion Questions

A If the electrical attraction between two pointlike objects at a distance of 1 m is 9×10^9 N, why can't we infer that their charges are +1 and -1 C? What further observations would we need to do in order to prove this?

B An electrically charged piece of tape will be attracted to your hand. Does that allow us to tell whether the mobile charged particles in your hand are positive or negative, or both?

8.1.3 Atoms

I was brought up to look at the atom as a nice, hard fellow, red or grey in color according to taste. *Rutherford*

Atomism

The Greeks have been kicked around a lot in the last couple of millennia: dominated by the Romans, bullied during the crusades by warlords going to and from the Holy Land, and occupied by Turkey until recently. It's no wonder they prefer to remember their salad days, when their best thinkers came up with concepts like democracy and atoms. Greece is democratic again after a period of military dictatorship, and an atom is proudly pictured on one of their coins. That's why it hurts me to have to say that the ancient Greek hypothesis that matter is made of atoms was pure guess-work. There was no real experimental evidence for atoms, and the 18th-century revival of the atom concept by Dalton owed little to the Greeks other than the name, which means "unsplittable." Subtracting even more cruelly from Greek glory, the name was shown to be inappropriate in 1897 when physicist J.J. Thomson proved experimentally that atoms had even smaller things inside them, which could be extracted. (Thomson called them "electrons.") The "unsplittable" was splittable after all.

But that's getting ahead of our story. What happened to the atom concept in the intervening two thousand years? Educated people continued to discuss the idea, and those who were in favor of it could often use it to give plausible explanations for various facts and phenomena. One fact that was readily explained was conservation of mass. For example, if you mix 1 kg of water with 1 kg of dirt, you get exactly 2 kg of mud, no more and no less. The same is true for a variety of processes such as freezing of water, fermenting beer, or pulverizing sandstone. If you believed in atoms, conservation of mass made perfect sense, because all these processes could be interpreted as mixing and rearranging atoms, without changing the total number of atoms. Still, this is nothing like a proof that atoms exist.

If atoms did exist, what types of atoms were there, and what dis-

tinguished the different types from each other? Was it their sizes, their shapes, their weights, or some other quality? The chasm between the ancient and modern atomisms becomes evident when we consider the wild speculations that existed on these issues until the present century. The ancients decided that there were four types of atoms, earth, water, air and fire; the most popular view was that they were distinguished by their shapes. Water atoms were spherical, hence water's ability to flow smoothly. Fire atoms had sharp points, which was why fire hurt when it touched one's skin. (There was no concept of temperature until thousands of years later.) The drastically different modern understanding of the structure of atoms was achieved in the course of the revolutionary decade stretching 1895 to 1905. The main purpose of this chapter is to describe those momentous experiments.

Atoms, light, and everything else

Although I tend to ridicule ancient Greek philosophers like Aristotle, let's take a moment to praise him for something. If you read Aristotle's writings on physics (or just skim them, which is all I've done), the most striking thing is how careful he is about classifying phenomena and analyzing relationships among phenomena. The human brain seems to naturally make a distinction between two types of physical phenomena: objects and motion of objects. When a phenomenon occurs that does not immediately present itself as one of these, there is a strong tendency to conceptualize it as one or the other, or even to ignore its existence completely. For instance, physics teachers shudder at students' statements that "the dynamite exploded, and force came out of it in all directions." In these examples, the nonmaterial concept of force is being mentally categorized as if it was a physical substance. The statement that "winding the clock stores motion in the spring" is a miscategorization of electrical energy as a form of motion. An example of ignoring the existence of a phenomenon altogether can be elicited by asking people why we need lamps. The typical response that "the lamp illuminates the room so we can see things," ignores the necessary role of light coming into our eyes from the things being illuminated.

If you ask someone to tell you briefly about atoms, the likely response is that "everything is made of atoms," but we've now seen that it's far from obvious which "everything" this statement would properly refer to. For the scientists of the early 1900s who were trying to investigate atoms, this was not a trivial issue of definitions. There was a new gizmo called the vacuum tube, of which the only familiar example today is the picture tube of a TV. In short order, electrical tinkerers had discovered a whole flock of new phenomena that occurred in and around vacuum tubes, and given them picturesque names like "x-rays," "cathode rays," "Hertzian waves," and "N-rays." These were the types of observations that ended up

telling us that we know about matter, but fierce controversies ensued over whether these were themselves forms of matter.

Let's bring ourselves up to the level of classification of phenomena employed by physicists in the year 1900. They recognized three categories:

- *Matter* has mass, can have kinetic energy, and can travel through a vacuum, transporting its mass and kinetic energy with it. Matter is conserved, both in the sense of conservation of mass and conservation of the number of atoms of each element. Atoms can't occupy the same space as other atoms, so a convenient way to prove something is not a form of matter is to show that it can pass through a solid material, in which the atoms are packed together closely.
- *Light* has no mass, always has energy, and can travel through a vacuum, transporting its energy with it. Two light beams can penetrate through each other and emerge from the collision without being weakened, deflected, or affected in any other way. Light can penetrate certain kinds of matter, e.g., glass.
- The third category is everything that doesn't fit the definition of light or matter. This catch-all category includes, for example, time, velocity, heat, and force.

The chemical elements

How would one find out what types of atoms there were? Today, it doesn't seem like it should have been very difficult to work out an experimental program to classify the types of atoms. For each type of atom, there should be a corresponding element, i.e., a pure substance made out of nothing but that type of atom. Atoms are supposed to be unsplittable, so a substance like milk could not possibly be elemental, since churning it vigorously causes it to split up into two separate substances: butter and whey. Similarly, rust could not be an element, because it can be made by combining two substances: iron and oxygen. Despite its apparent reasonableness, no such program was carried out until the eighteenth century. The ancients presumably did not do it because observation was not universally agreed on as the right way to answer questions about nature, and also because they lacked the necessary techniques or the techniques were the province of laborers with low social status, such as smiths and miners. Alchemists were hindered by atomism's reputation for subversiveness, and by a tendency toward mysticism and secrecy. (The most celebrated challenge facing the alchemists, that of converting lead into gold, is one we now know to be impossible, since lead and gold are both elements.)

$$\frac{m_{\text{He}}}{m_{\text{H}}} = 3.97$$

$$\frac{m_{\text{Ne}}}{m_{\text{H}}} = 20.01$$

$$\frac{m_{\text{Sc}}}{m_{\text{H}}} = 44.60$$

d / Examples of masses of atoms compared to that of hydrogen. Note how some, but not all, are close to integers.

By 1900, however, chemists had done a reasonably good job of finding out what the elements were. They also had determined the

ratios of the different atoms' masses fairly accurately. A typical technique would be to measure how many grams of sodium (Na) would combine with one gram of chlorine (Cl) to make salt (NaCl). (This assumes you've already decided based on other evidence that salt consisted of equal numbers of Na and Cl atoms.) The masses of individual atoms, as opposed to the mass ratios, were known only to within a few orders of magnitude based on indirect evidence, and plenty of physicists and chemists denied that individual atoms were anything more than convenient symbols.

Making sense of the elements

As the information accumulated, the challenge was to find a way of systematizing it; the modern scientist's aesthetic sense rebels against complication. This hodgepodge of elements was an embarrassment. One contemporary observer, William Crookes, described the elements as extending "before us as stretched the wide Atlantic before the gaze of Columbus, mocking, taunting and murmuring strange riddles, which no man has yet been able to solve." It wasn't long before people started recognizing that many atoms' masses were nearly integer multiples of the mass of hydrogen, the lightest element. A few excitable types began speculating that hydrogen was the basic building block, and that the heavier elements were made of clusters of hydrogen. It wasn't long, however, before their parade was rained on by more accurate measurements, which showed that not all of the elements had atomic masses that were near integer multiples of hydrogen, and even the ones that were close to being integer multiples were off by one percent or so.

¹ H																			² He	
³ Li	⁴ Be																			
¹¹ Na	¹² Mg																			
¹⁹ K	²⁰ Ca	²¹ Sc	²² Ti	²³ V	²⁴ Cr	²⁵ Mn	²⁶ Fe	²⁷ Co	²⁸ Ni	²⁹ Cu	³⁰ Zn	³¹ Ga	³² Ge	³³ As	³⁴ Se	³⁵ Br	³⁶ Kr			
³⁷ Rb	³⁸ Sr	³⁹ Y	⁴⁰ Zr	⁴¹ Nb	⁴² Mo	⁴³ Tc	⁴⁴ Ru	⁴⁵ Rh	⁴⁶ Pd	⁴⁷ Ag	⁴⁸ Cd	⁴⁹ In	⁵⁰ Sn	⁵¹ Sb	⁵² Te	⁵³ I	⁵⁴ Xe			
⁵⁵ Cs	⁵⁶ Ba	⁵⁷ La *	⁷² Hf	⁷³ Ta	⁷⁴ W	⁷⁵ Re	⁷⁶ Os	⁷⁷ Ir	⁷⁸ Pt	⁷⁹ Au	⁸⁰ Hg	⁸¹ Tl	⁸² Pb	⁸³ Bi	⁸⁴ Po	⁸⁵ At	⁸⁶ Rn			
⁸⁷ Fr	⁸⁸ Ra	⁸⁹ Ac **	¹⁰⁴ Rf	¹⁰⁵ Ha	¹⁰⁶	¹⁰⁷	¹⁰⁸	¹⁰⁹	¹¹⁰	¹¹¹	¹¹²	¹¹³	¹¹⁴	¹¹⁵	¹¹⁶	¹¹⁷	¹¹⁸			
*	⁵⁸ Ce	⁵⁹ Pr	⁶⁰ Nd	⁶¹ Pm	⁶² Sm	⁶³ Eu	⁶⁴ Gd	⁶⁵ Tb	⁶⁶ Dy	⁶⁷ Ho	⁶⁸ Er	⁶⁹ Tm	⁷⁰ Yb	⁷¹ Lu						
**	⁹⁰ Th	⁹¹ Pa	⁹² U	⁹³ Np	⁹⁴ Pu	⁹⁵ Am	⁹⁶ Cm	⁹⁷ Bk	⁹⁸ Cf	⁹⁹ Es	¹⁰⁰ Fm	¹⁰¹ Md	¹⁰² No	¹⁰³ Lr						

Chemistry professor Dmitri Mendeleev, preparing his lectures in 1869, wanted to find some way to organize his knowledge for his students to make it more understandable. He wrote the names of all the elements on cards and began arranging them in different ways on his desk, trying to find an arrangement that would make sense of

e / A modern periodic table. Elements in the same column have similar chemical properties. The modern atomic numbers, discussed in section 8.2, were not known in Mendeleev's time, since the table could be flipped in various ways.

the muddle. The row-and-column scheme he came up with is essentially our modern periodic table. The columns of the modern version represent groups of elements with similar chemical properties, and each row is more massive than the one above it. Going across each row, this almost always resulted in placing the atoms in sequence by weight as well. What made the system significant was its predictive value. There were three places where Mendeleev had to leave gaps in his checkerboard to keep chemically similar elements in the same column. He predicted that elements would exist to fill these gaps, and extrapolated or interpolated from other elements in the same column to predict their numerical properties, such as masses, boiling points, and densities. Mendeleev's professional stock skyrocketed when his three elements (later named gallium, scandium and germanium) were discovered and found to have very nearly the properties he had predicted.

One thing that Mendeleev's table made clear was that mass was not the basic property that distinguished atoms of different elements. To make his table work, he had to deviate from ordering the elements strictly by mass. For instance, iodine atoms are lighter than tellurium, but Mendeleev had to put iodine after tellurium so that it would lie in a column with chemically similar elements.

Direct proof that atoms existed

The success of the kinetic theory of heat was taken as strong evidence that, in addition to the motion of any object as a whole, there is an invisible type of motion all around us: the random motion of atoms within each object. But many conservatives were not convinced that atoms really existed. Nobody had ever seen one, after all. It wasn't until generations after the kinetic theory of heat was developed that it was demonstrated conclusively that atoms really existed and that they participated in continuous motion that never died out.

The smoking gun to prove atoms were more than mathematical abstractions came when some old, obscure observations were reexamined by an unknown Swiss patent clerk named Albert Einstein. A botanist named Brown, using a microscope that was state of the art in 1827, observed tiny grains of pollen in a drop of water on a microscope slide, and found that they jumped around randomly for no apparent reason. Wondering at first if the pollen he'd assumed to be dead was actually alive, he tried looking at particles of soot, and found that the soot particles also moved around. The same results would occur with any small grain or particle suspended in a liquid. The phenomenon came to be referred to as Brownian motion, and its existence was filed away as a quaint and thoroughly unimportant fact, really just a nuisance for the microscopist.

It wasn't until 1906 that Einstein found the correct interpreta-

tion for Brown's observation: the water molecules were in continuous random motion, and were colliding with the particle all the time, kicking it in random directions. After all the millennia of speculation about atoms, at last there was solid proof. Einstein's calculations dispelled all doubt, since he was able to make accurate predictions of things like the average distance traveled by the particle in a certain amount of time. (Einstein received the Nobel Prize not for his theory of relativity but for his papers on Brownian motion and the photoelectric effect.)

Discussion Questions

- A How could knowledge of the size of an individual aluminum atom be used to infer an estimate of its mass, or vice versa?
- B How could one test Einstein's interpretation of Brownian motion by observing it at different temperatures?

8.1.4 Quantization of charge

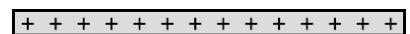
Proving that atoms actually existed was a big accomplishment, but demonstrating their existence was different from understanding their properties. Note that the Brown-Einstein observations had nothing at all to do with electricity, and yet we know that matter is inherently electrical, and we have been successful in interpreting certain electrical phenomena in terms of mobile positively and negatively charged particles. Are these particles atoms? Parts of atoms? Particles that are entirely separate from atoms? It is perhaps premature to attempt to answer these questions without any conclusive evidence in favor of the charged-particle model of electricity.

Strong support for the charged-particle model came from a 1911 experiment by physicist Robert Millikan at the University of Chicago. Consider a jet of droplets of perfume or some other liquid made by blowing it through a tiny pinhole. The droplets emerging from the pinhole must be smaller than the pinhole, and in fact most of them are even more microscopic than that, since the turbulent flow of air tends to break them up. Millikan reasoned that the droplets would acquire a little bit of electric charge as they rubbed against the channel through which they emerged, and if the charged-particle model of electricity was right, the charge might be split up among so many minuscule liquid drops that a single drop might have a total charge amounting to an excess of only a few charged particles — perhaps an excess of one positive particle on a certain drop, or an excess of two negative ones on another.

Millikan's ingenious apparatus, g, consisted of two metal plates, which could be electrically charged as needed. He sprayed a cloud of oil droplets into the space between the plates, and selected one drop through a microscope for study. First, with no charge on the plates, he would determine the drop's mass by letting it fall through the air and measuring its terminal velocity, i.e., the velocity at which



f / A young Robert Millikan. (*Contemporary*)



g / A simplified diagram of Millikan's apparatus.

the force of air friction canceled out the force of gravity. The force of air drag on a slowly moving sphere had already been found by experiment to be bvr^2 , where b was a constant. Setting the total force equal to zero when the drop is at terminal velocity gives

$$bvr^2 - mg = 0,$$

and setting the known density of oil equal to the drop's mass divided by its volume gives a second equation,

$$\rho = \frac{m}{\frac{4}{3}\pi r^3}.$$

Everything in these equations can be measured directly except for m and r , so these are two equations in two unknowns, which can be solved in order to determine how big the drop is.

Next Millikan charged the metal plates, adjusting the amount of charge so as to exactly counteract gravity and levitate the drop. If, for instance, the drop being examined happened to have a total charge that was negative, then positive charge put on the top plate would attract it, pulling it up, and negative charge on the bottom plate would repel it, pushing it up. (Theoretically only one plate would be necessary, but in practice a two-plate arrangement like this gave electrical forces that were more uniform in strength throughout the space where the oil drops were.) The amount of charge on the plates required to levitate the charged drop gave Millikan a handle on the amount of charge the drop carried. The more charge the drop had, the stronger the electrical forces on it would be, and the less charge would have to be put on the plates to do the trick. Unfortunately, expressing this relationship using Coulomb's law would have been impractical, because it would require a perfect knowledge of how the charge was distributed on each plate, plus the ability to perform vector addition of all the forces being exerted on the drop by all the charges on the plate. Instead, Millikan made use of the fact that the electrical force experienced by a pointlike charged object at a certain point in space is proportional to its charge,

$$\frac{F}{q} = \text{constant.}$$

q	$/(1.64 \times 10^{-19} \text{ C})$
-1.970×10^{-18}	-12.02
-0.987×10^{-18}	-6.02
-2.773×10^{-18}	-16.93

h / A few samples of Millikan's data.

With a given amount of charge on the plates, this constant could be determined for instance by discarding the oil drop, inserting between the plates a larger and more easily handled object with a known charge on it, and measuring the force with conventional methods. (Millikan actually used a slightly different set of techniques for determining the constant, but the concept is the same.) The amount of force on the actual oil drop had to equal mg , since it was just enough to levitate it, and once the calibration constant had been determined, the charge of the drop could then be found based on its previously determined mass.

The table on the left shows a few of the results from Millikan's 1911 paper. (Millikan took data on both negatively and positively charged drops, but in his paper he gave only a sample of his data on negatively charged drops, so these numbers are all negative.) Even a quick look at the data leads to the suspicion that the charges are not simply a series of random numbers. For instance, the second charge is almost exactly equal to half the first one. Millikan explained the observed charges as all being integer multiples of a single number, 1.64×10^{-19} C. In the second column, dividing by this constant gives numbers that are essentially integers, allowing for the random errors present in the experiment. Millikan states in his paper that these results were a

...direct and tangible demonstration ... of the correctness of the view advanced many years ago and supported by evidence from many sources that all electrical charges, however produced, are exact multiples of one definite, elementary electrical charge, or in other words, that an electrical charge instead of being spread uniformly over the charged surface has a definite granular structure, consisting, in fact, of ... specks, or atoms of electricity, all precisely alike, peppered over the surface of the charged body.

In other words, he had provided direct evidence for the charged-particle model of electricity and against models in which electricity was described as some sort of fluid. The basic charge is notated e , and the modern value is $e = 1.60 \times 10^{-19}$ C. The word “quantized” is used in physics to describe a quantity that can only have certain numerical values, and cannot have any of the values between those. In this language, we would say that Millikan discovered that charge is quantized. The charge e is referred to as the quantum of charge.

A historical note on Millikan's fraud

Very few undergraduate physics textbooks mention the well-documented fact that although Millikan's conclusions were correct, he was guilty of scientific fraud. His technique was difficult and painstaking to perform, and his original notebooks, which have been preserved, show that the data were far less perfect than he claimed in his published scientific papers. In his publications, he stated categorically that every single oil drop observed had had a charge that was a multiple of e , with no exceptions or omissions. But his notebooks are replete with notations such as “beautiful data, keep,” and “bad run, throw out.” Millikan, then, appears to have earned his Nobel Prize by advocating a correct position with dishonest descriptions of his data.

Why do textbook authors fail to mention Millikan's fraud? It may be that they think students are too unsophisticated to cor-

rectly evaluate the implications of the fact that scientific fraud has sometimes existed and even been rewarded by the scientific establishment. Maybe they are afraid students will reason that fudging data is OK, since Millikan got the Nobel Prize for it. But falsifying history in the name of encouraging truthfulness is more than a little ironic. English teachers don't edit Shakespeare's tragedies so that the bad characters are always punished and the good ones never suffer!

self-check C

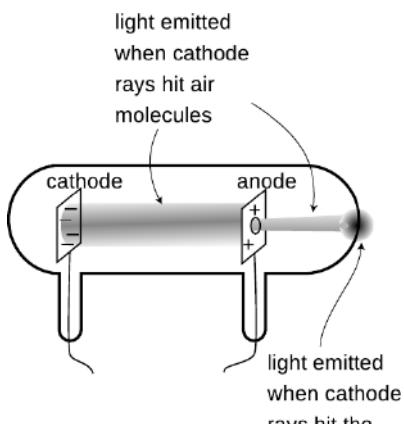
Is money quantized? What is the quantum of money? ▷ Answer, p. 1062

8.1.5 The electron

Cathode rays

Nineteenth-century physicists spent a lot of time trying to come up with wild, random ways to play with electricity. The best experiments of this kind were the ones that made big sparks or pretty colors of light.

One such parlor trick was the cathode ray. To produce it, you first had to hire a good glassblower and find a good vacuum pump. The glassblower would create a hollow tube and embed two pieces of metal in it, called the electrodes, which were connected to the outside via metal wires passing through the glass. Before letting him seal up the whole tube, you would hook it up to a vacuum pump, and spend several hours huffing and puffing away at the pump's hand crank to get a good vacuum inside. Then, while you were still pumping on the tube, the glassblower would melt the glass and seal the whole thing shut. Finally, you would put a large amount of positive charge on one wire and a large amount of negative charge on the other. Metals have the property of letting charge move through them easily, so the charge deposited on one of the wires would quickly spread out because of the repulsion of each part of it for every other part. This spreading-out process would result in nearly all the charge ending up in the electrodes, where there is more room to spread out than there is in the wire. For obscure historical reasons a negative electrode is called a cathode and a positive one is an anode.



i / Cathode rays observed in a vacuum tube.

Figure i shows the light-emitting stream that was observed. If, as shown in this figure, a hole was made in the anode, the beam would extend on through the hole until it hit the glass. Drilling a hole in the cathode, however would not result in any beam coming out on the left side, and this indicated that the stuff, whatever it was, was coming from the cathode. The rays were therefore christened "cathode rays." (The terminology is still used today in the term "cathode ray tube" or "CRT" for the picture tube of a TV or computer monitor.)

Were cathode rays a form of light, or of matter?

Were cathode rays a form of light, or matter? At first no one really cared what they were, but as their scientific importance became more apparent, the light-versus-matter issue turned into a controversy along nationalistic lines, with the Germans advocating light and the English holding out for matter. The supporters of the material interpretation imagined the rays as consisting of a stream of atoms ripped from the substance of the cathode.

One of our defining characteristics of matter is that material objects cannot pass through each other. Experiments showed that cathode rays could penetrate at least some small thickness of matter, such as a metal foil a tenth of a millimeter thick, implying that they were a form of light.

Other experiments, however, pointed to the contrary conclusion. Light is a wave phenomenon, and one distinguishing property of waves is demonstrated by speaking into one end of a paper towel roll. The sound waves do not emerge from the other end of the tube as a focused beam. Instead, they begin spreading out in all directions as soon as they emerge. This shows that waves do not necessarily travel in straight lines. If a piece of metal foil in the shape of a star or a cross was placed in the way of the cathode ray, then a “shadow” of the same shape would appear on the glass, showing that the rays traveled in straight lines. This straight-line motion suggested that they were a stream of small particles of matter.

These observations were inconclusive, so what was really needed was a determination of whether the rays had mass and weight. The trouble was that cathode rays could not simply be collected in a cup and put on a scale. When the cathode ray tube is in operation, one does not observe any loss of material from the cathode, or any crust being deposited on the anode.

Nobody could think of a good way to weigh cathode rays, so the next most obvious way of settling the light/matter debate was to check whether the cathode rays possessed electrical charge. Light was known to be uncharged. If the cathode rays carried charge, they were definitely matter and not light, and they were presumably being made to jump the gap by the simultaneous repulsion of the negative charge in the cathode and attraction of the positive charge in the anode. The rays would overshoot the anode because of their momentum. (Although electrically charged particles do not normally leap across a gap of vacuum, very large amounts of charge were being used, so the forces were unusually intense.)

Thomson's experiments

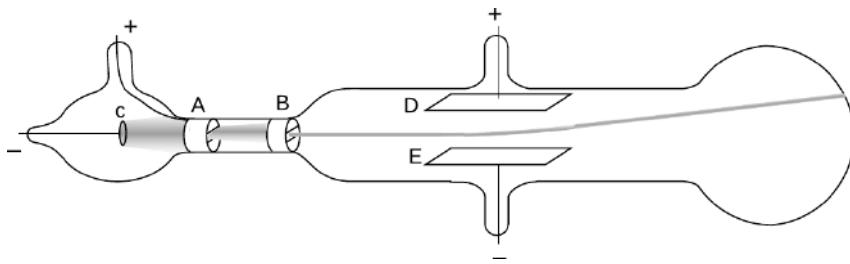
Physicist J.J. Thomson at Cambridge carried out a series of definitive experiments on cathode rays around the year 1897. By turning them slightly off course with electrical forces, k, he showed



j / J.J. Thomson in the lab.

that they were indeed electrically charged, which was strong evidence that they were material. Not only that, but he proved that they had mass, and measured the ratio of their mass to their charge, m/q . Since their mass was not zero, he concluded that they were a form of matter, and presumably made up of a stream of microscopic, negatively charged particles. When Millikan published his results fourteen years later, it was reasonable to assume that the charge of one such particle equaled minus one fundamental charge, $q = -e$, and from the combination of Thomson's and Millikan's results one could therefore determine the mass of a single cathode ray particle.

k / Thomson's experiment proving cathode rays had electric charge (redrawn from his original paper). The cathode, C, and anode, A, are as in any cathode ray tube. The rays pass through a slit in the anode, and a second slit, B, is interposed in order to make the beam thinner and eliminate rays that were not going straight. Charging plates D and E shows that cathode rays have charge: they are attracted toward the positive plate D and repelled by the negative plate E.



The basic technique for determining m/q was simply to measure the angle through which the charged plates bent the beam. The electric force acting on a cathode ray particle while it was between the plates would be proportional to its charge,

$$F_{elec} = (\text{known constant}) \cdot q.$$

Application of Newton's second law, $a = F/m$, would allow m/q to be determined:

$$\frac{m}{q} = \frac{\text{known constant}}{a}$$

There was just one catch. Thomson needed to know the cathode ray particles' velocity in order to figure out their acceleration. At that point, however, nobody had even an educated guess as to the speed of the cathode rays produced in a given vacuum tube. The beam appeared to leap across the vacuum tube practically instantaneously, so it was no simple matter of timing it with a stopwatch!

Thomson's clever solution was to observe the effect of both electric and magnetic forces on the beam. The magnetic force exerted by a particular magnet would depend on both the cathode ray's charge and its velocity:

$$F_{mag} = (\text{known constant } \#2) \cdot qv$$

Thomson played with the electric and magnetic forces until either one would produce an equal effect on the beam, allowing him to solve for the velocity,

$$v = \frac{\text{(known constant)}}{\text{(known constant \#2)}}.$$

Knowing the velocity (which was on the order of 10% of the speed of light for his setup), he was able to find the acceleration and thus the mass-to-charge ratio m/q . Thomson's techniques were relatively crude (or perhaps more charitably we could say that they stretched the state of the art of the time), so with various methods he came up with m/q values that ranged over about a factor of two, even for cathode rays extracted from a cathode made of a single material. The best modern value is $m/q = 5.69 \times 10^{-12}$ kg/C, which is consistent with the low end of Thomson's range.

The cathode ray as a subatomic particle: the electron

What was significant about Thomson's experiment was not the actual numerical value of m/q , however, so much as the fact that, combined with Millikan's value of the fundamental charge, it gave a mass for the cathode ray particles that was thousands of times smaller than the mass of even the lightest atoms. Even without Millikan's results, which were 14 years in the future, Thomson recognized that the cathode rays' m/q was thousands of times smaller than the m/q ratios that had been measured for electrically charged atoms in chemical solutions. He correctly interpreted this as evidence that the cathode rays were smaller building blocks — he called them *electrons* — out of which atoms themselves were formed. This was an extremely radical claim, coming at a time when atoms had not yet been proven to exist! Even those who used the word “atom” often considered them no more than mathematical abstractions, not literal objects. The idea of searching for structure inside of “un-splittable” atoms was seen by some as lunacy, but within ten years Thomson's ideas had been amply verified by many more detailed experiments.

Discussion Questions

A Thomson started to become convinced during his experiments that the “cathode rays” observed coming from the cathodes of vacuum tubes were building blocks of atoms — what we now call electrons. He then carried out observations with cathodes made of a variety of metals, and found that m/q was roughly the same in every case, considering his limited accuracy. Given his suspicion, why did it make sense to try different metals? How would the consistent values of m/q serve to test his hypothesis?

B My students have frequently asked whether the m/q that Thomson measured was the value for a single electron, or for the whole beam. Can you answer this question?

C Thomson found that the m/q of an electron was thousands of times smaller than that of charged atoms in chemical solutions. Would this imply that the electrons had more charge? Less mass? Would there be no way to tell? Explain. Remember that Millikan's results were still many years in the future, so q was unknown.

D Can you guess any practical reason why Thomson couldn't just let one electron fly across the gap before disconnecting the battery and turning off the beam, and then measure the amount of charge deposited on the anode, thus allowing him to measure the charge of a single electron directly?

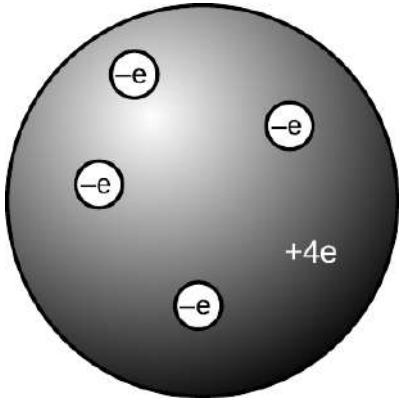
E Why is it not possible to determine m and q themselves, rather than just their ratio, by observing electrons' motion in electric and magnetic fields?

8.1.6 The raisin cookie model of the atom

Based on his experiments, Thomson proposed a picture of the atom which became known as the raisin cookie model. In the neutral atom, l, there are four electrons with a total charge of $-4e$, sitting in a sphere (the "cookie") with a charge of $+4e$ spread throughout it. It was known that chemical reactions could not change one element into another, so in Thomson's scenario, each element's cookie sphere had a permanently fixed radius, mass, and positive charge, different from those of other elements. The electrons, however, were not a permanent feature of the atom, and could be tacked on or pulled out to make charged ions. Although we now know, for instance, that a neutral atom with four electrons is the element beryllium, scientists at the time did not know how many electrons the various neutral atoms possessed.

This model is clearly different from the one you've learned in grade school or through popular culture, where the positive charge is concentrated in a tiny nucleus at the atom's center. An equally important change in ideas about the atom has been the realization that atoms and their constituent subatomic particles behave entirely differently from objects on the human scale. For instance, we'll see later that an electron can be in more than one place at one time. The raisin cookie model was part of a long tradition of attempts to make mechanical models of phenomena, and Thomson and his contemporaries never questioned the appropriateness of building a mental model of an atom as a machine with little parts inside. Today, mechanical models of atoms are still used (for instance the tinker-toy-style molecular modeling kits like the ones used by Watson and Crick to figure out the double helix structure of DNA), but scientists realize that the physical objects are only aids to help our brains' symbolic and visual processes think about atoms.

Although there was no clear-cut experimental evidence for many of the details of the raisin cookie model, physicists went ahead and started working out its implications. For instance, suppose you had



I / The raisin cookie model of the atom with four units of charge, which we now know to be beryllium.

a four-electron atom. All four electrons would be repelling each other, but they would also all be attracted toward the center of the “cookie” sphere. The result should be some kind of stable, symmetric arrangement in which all the forces canceled out. People sufficiently clever with math soon showed that the electrons in a four-electron atom should settle down at the vertices of a pyramid with one less side than the Egyptian kind, i.e., a regular tetrahedron. This deduction turns out to be wrong because it was based on incorrect features of the model, but the model also had many successes, a few of which we will now discuss.

Flow of electrical charge in wires

example 3

One of my former students was the son of an electrician, and had become an electrician himself. He related to me how his father had remained refused to believe all his life that electrons really flowed through wires. If they had, he reasoned, the metal would have gradually become more and more damaged, eventually crumbling to dust.

His opinion is not at all unreasonable based on the fact that electrons are material particles, and that matter cannot normally pass through matter without making a hole through it. Nineteenth-century physicists would have shared his objection to a charged-particle model of the flow of electrical charge. In the raisin-cookie model, however, the electrons are very low in mass, and therefore presumably very small in size as well. It is not surprising that they can slip between the atoms without damaging them.

Flow of electrical charge across cell membranes

example 4

Your nervous system is based on signals carried by charge moving from nerve cell to nerve cell. Your body is essentially all liquid, and atoms in a liquid are mobile. This means that, unlike the case of charge flowing in a solid wire, entire charged atoms can flow in your nervous system

Emission of electrons in a cathode ray tube

example 5

Why do electrons detach themselves from the cathode of a vacuum tube? Certainly they are encouraged to do so by the repulsion of the negative charge placed on the cathode and the attraction from the net positive charge of the anode, but these are not strong enough to rip electrons out of atoms by main force — if they were, then the entire apparatus would have been instantly vaporized as every atom was simultaneously ripped apart!

The raisin cookie model leads to a simple explanation. We know that heat is the energy of random motion of atoms. The atoms in any object are therefore violently jostling each other all the time, and a few of these collisions are violent enough to knock electrons out of atoms. If this occurs near the surface of a solid object, the electron may come loose. Ordinarily, however, this loss of electrons is a self-limiting process; the loss of electrons leaves

the object with a net positive charge, which attracts the lost sheep home to the fold. (For objects immersed in air rather than vacuum, there will also be a balanced exchange of electrons between the air and the object.)

This interpretation explains the warm and friendly yellow glow of the vacuum tubes in an antique radio. To encourage the emission of electrons from the vacuum tubes' cathodes, the cathodes are intentionally warmed up with little heater coils.

Discussion Questions

A Today many people would define an ion as an atom (or molecule) with missing electrons or extra electrons added on. How would people have defined the word "ion" before the discovery of the electron?

B Since electrically neutral atoms were known to exist, there had to be positively charged subatomic stuff to cancel out the negatively charged electrons in an atom. Based on the state of knowledge immediately after the Millikan and Thomson experiments, was it possible that the positively charged stuff had an unquantized amount of charge? Could it be quantized in units of $+e$? In units of $+2e$? In units of $+5/7e$?

This chapter is summarized on page 1084. Notation and terminology are tabulated on pages 1070-1071.

8.2 The nucleus

8.2.1 Radioactivity

Becquerel's discovery of radioactivity

How did physicists figure out that the raisin cookie model was incorrect, and that the atom's positive charge was concentrated in a tiny, central nucleus? The story begins with the discovery of radioactivity by the French chemist Becquerel. Up until radioactivity was discovered, all the processes of nature were thought to be based on chemical reactions, which were rearrangements of combinations of atoms. Atoms exert forces on each other when they are close together, so sticking or unsticking them would either release or store electrical energy. That energy could be converted to and from other forms, as when a plant uses the energy in sunlight to make sugars and carbohydrates, or when a child eats sugar, releasing the energy in the form of kinetic energy.

Becquerel discovered a process that seemed to release energy from an unknown new source that was not chemical. Becquerel, whose father and grandfather had also been physicists, spent the first twenty years of his professional life as a successful civil engineer, teaching physics on a part-time basis. He was awarded the chair of physics at the Musée d'Histoire Naturelle in Paris after the death of his father, who had previously occupied it. Having now a significant amount of time to devote to physics, he began studying the interaction of light and matter. He became interested in the phe-

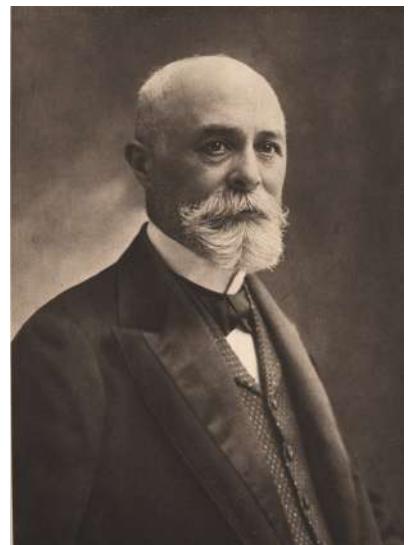
nomenon of phosphorescence, in which a substance absorbs energy from light, then releases the energy via a glow that only gradually goes away. One of the substances he investigated was a uranium compound, the salt UOSO_5 . One day in 1896, cloudy weather interfered with his plan to expose this substance to sunlight in order to observe its fluorescence. He stuck it in a drawer, coincidentally on top of a blank photographic plate — the old-fashioned glass-backed counterpart of the modern plastic roll of film. The plate had been carefully wrapped, but several days later when Becquerel checked it in the darkroom before using it, he found that it was ruined, as if it had been completely exposed to light.

History provides many examples of scientific discoveries that happened this way: an alert and inquisitive mind decides to investigate a phenomenon that most people would not have worried about explaining. Becquerel first determined by further experiments that the effect was produced by the uranium salt, despite a thick wrapping of paper around the plate that blocked out all light. He tried a variety of compounds, and found that it was the uranium that did it: the effect was produced by any uranium compound, but not by any compound that didn't include uranium atoms. The effect could be at least partially blocked by a sufficient thickness of metal, and he was able to produce silhouettes of coins by interposing them between the uranium and the plate. This indicated that the effect traveled in a straight line., so that it must have been some kind of ray rather than, e.g., the seepage of chemicals through the paper. He used the word “radiations,” since the effect radiated out from the uranium salt.

At this point Becquerel still believed that the uranium atoms were absorbing energy from light and then gradually releasing the energy in the form of the mysterious rays, and this was how he presented it in his first published lecture describing his experiments. Interesting, but not earth-shattering. But he then tried to determine how long it took for the uranium to use up all the energy that had supposedly been stored in it by light, and he found that it never seemed to become inactive, no matter how long he waited. Not only that, but a sample that had been exposed to intense sunlight for a whole afternoon was no more or less effective than a sample that had always been kept inside. Was this a violation of conservation of energy? If the energy didn't come from exposure to light, where did it come from?

Three kinds of “radiations”

Unable to determine the source of the energy directly, turn-of-the-century physicists instead studied the behavior of the “radiations” once they had been emitted. Becquerel had already shown that the radioactivity could penetrate through cloth and paper, so the first obvious thing to do was to investigate in more detail what



a / Henri Becquerel (1852-1908).



b / Becquerel's photographic plate. In the exposure at the bottom of the image, he has found that he could absorb the radiations, casting the shadow of a Maltese cross that was placed between the plate and the uranium salts.

thickness of material the radioactivity could get through. They soon learned that a certain fraction of the radioactivity's intensity would be eliminated by even a few inches of air, but the remainder was not eliminated by passing through more air. Apparently, then, the radioactivity was a mixture of more than one type, of which one was blocked by air. They then found that of the part that could penetrate air, a further fraction could be eliminated by a piece of paper or a very thin metal foil. What was left after that, however, was a third, extremely penetrating type, some of whose intensity would still remain even after passing through a brick wall. They decided that this showed there were three types of radioactivity, and without having the faintest idea of what they really were, they made up names for them. The least penetrating type was arbitrarily labeled α (alpha), the first letter of the Greek alphabet, and so on through β (beta) and finally γ (gamma) for the most penetrating type.

Radium: a more intense source of radioactivity

The measuring devices used to detect radioactivity were crude: photographic plates or even human eyeballs (radioactivity makes flashes of light in the jelly-like fluid inside the eye, which can be seen by the eyeball's owner if it is otherwise very dark). Because the ways of detecting radioactivity were so crude and insensitive, further progress was hindered by the fact that the amount of radioactivity emitted by uranium was not really very great. The vital contribution of physicist/chemist Marie Curie and her husband Pierre was to discover the element radium, and to purify and isolate significant quantities of it. Radium emits about a million times more radioactivity per unit mass than uranium, making it possible to do the experiments that were needed to learn the true nature of radioactivity. The dangers of radioactivity to human health were then unknown, and Marie died of leukemia thirty years later. (Pierre was run over and killed by a horsecart.)

Tracking down the nature of alphas, betas, and gammas

As radium was becoming available, an apprentice scientist named Ernest Rutherford arrived in England from his native New Zealand and began studying radioactivity at the Cavendish Laboratory. The young colonial's first success was to measure the mass-to-charge ratio of beta rays. The technique was essentially the same as the one Thomson had used to measure the mass-to-charge ratio of cathode rays by measuring their deflections in electric and magnetic fields. The only difference was that instead of the cathode of a vacuum tube, a nugget of radium was used to supply the beta rays. Not only was the technique the same, but so was the result. Beta rays had the same m/q ratio as cathode rays, which suggested they were one and the same. Nowadays, it would make sense simply to use the term "electron," and avoid the archaic "cathode ray" and "beta particle," but the old labels are still widely used, and it is unfortu-

nately necessary for physics students to memorize all three names for the same thing.

At first, it seemed that neither alphas or gammas could be deflected in electric or magnetic fields, making it appear that neither was electrically charged. But soon Rutherford obtained a much more powerful magnet, and was able to use it to deflect the alphas but not the gammas. The alphas had a much larger value of m/q than the betas (about 4000 times greater), which was why they had been so hard to deflect. Gammas are uncharged, and were later found to be a form of light.

The m/q ratio of alpha particles turned out to be the same as those of two different types of ions, He^{++} (a helium atom with two missing electrons) and H_2^+ (two hydrogen atoms bonded into a molecule, with one electron missing), so it seemed likely that they were one or the other of those. The diagram shows a simplified version of Rutherford's ingenious experiment proving that they were He^{++} ions. The gaseous element radon, an alpha emitter, was introduced into one half of a double glass chamber. The glass wall dividing the chamber was made extremely thin, so that some of the rapidly moving alpha particles were able to penetrate it. The other chamber, which was initially evacuated, gradually began to accumulate a population of alpha particles (which would quickly pick up electrons from their surroundings and become electrically neutral). Rutherford then determined that it was helium gas that had appeared in the second chamber. Thus alpha particles were proved to be He^{++} ions. The nucleus was yet to be discovered, but in modern terms, we would describe a He^{++} ion as the nucleus of a He atom.

To summarize, here are the three types of radiation emitted by radioactive elements, and their descriptions in modern terms:

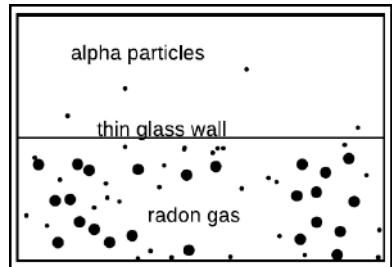
α particle	stopped by a few inches of air	He nucleus
β particle	stopped by a piece of paper	electron
γ ray	penetrates thick shielding	a type of light

Discussion Question

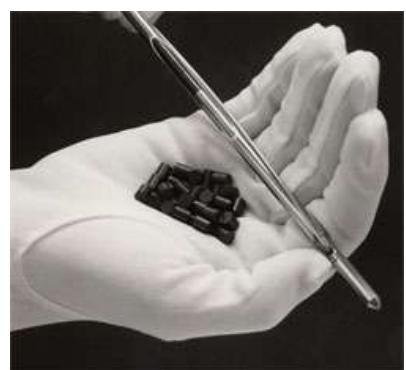
A Most sources of radioactivity emit alphas, betas, and gammas, not just one of the three. In the radon experiment, how did Rutherford know that he was studying the alphas?

8.2.2 The planetary model

The stage was now set for the unexpected discovery that the positively charged part of the atom was a tiny, dense lump at the atom's center rather than the "cookie dough" of the raisin cookie model. By 1909, Rutherford was an established professor, and had students working under him. For a raw undergraduate named Marsden, he picked a research project he thought would be tedious but straightforward.



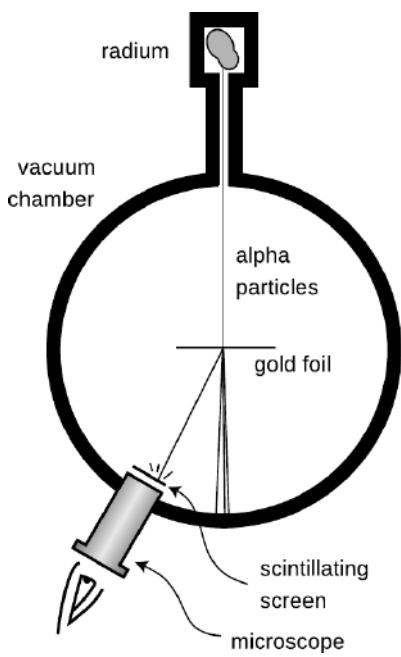
c / A simplified version of Rutherford's 1908 experiment, showing that alpha particles were doubly ionized helium atoms.



d / These pellets of uranium fuel will be inserted into the metal fuel rod and used in a nuclear reactor. The pellets emit alpha and beta radiation, which the gloves are thick enough to stop.



e / Ernest Rutherford (1871-1937).



f / Marsden and Rutherford's apparatus.

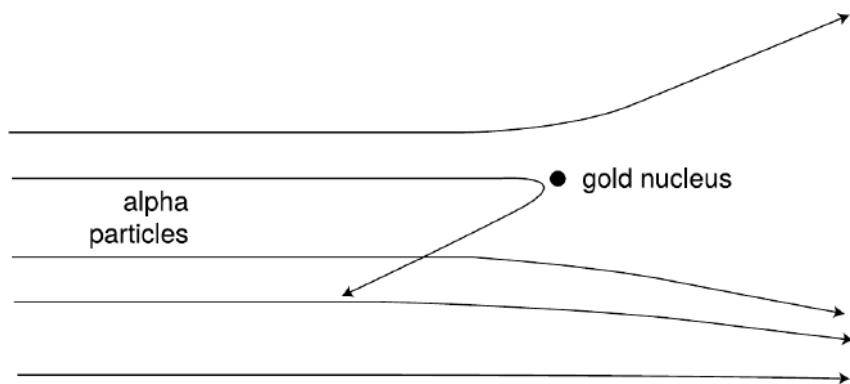
It was already known that although alpha particles would be stopped completely by a sheet of paper, they could pass through a sufficiently thin metal foil. Marsden was to work with a gold foil only 1000 atoms thick. (The foil was probably made by evaporating a little gold in a vacuum chamber so that a thin layer would be deposited on a glass microscope slide. The foil would then be lifted off the slide by submerging the slide in water.)

Rutherford had already determined in his previous experiments the speed of the alpha particles emitted by radium, a fantastic 1.5×10^7 m/s. The experimenters in Rutherford's group visualized them as very small, very fast cannonballs penetrating the "cookie dough" part of the big gold atoms. A piece of paper has a thickness of a hundred thousand atoms or so, which would be sufficient to stop them completely, but crashing through a thousand would only slow them a little and turn them slightly off of their original paths.

Marsden's supposedly ho-hum assignment was to use the apparatus shown in figure f to measure how often alpha particles were deflected at various angles. A tiny lump of radium in a box emitted alpha particles, and a thin beam was created by blocking all the alphas except those that happened to pass out through a tube. Typically deflected in the gold by only a small amount, they would reach a screen very much like the screen of a TV's picture tube, which would make a flash of light when it was hit. Here is the first example we have encountered of an experiment in which a beam of particles is detected one at a time. This was possible because each alpha particle carried so much kinetic energy; they were moving at about the same speed as the electrons in the Thomson experiment, but had ten thousand times more mass.

Marsden sat in a dark room, watching the apparatus hour after hour and recording the number of flashes with the screen moved to various angles. The rate of the flashes was highest when he set the screen at an angle close to the line of the alphas' original path, but if he watched an area farther off to the side, he would also occasionally see an alpha that had been deflected through a larger angle. After seeing a few of these, he got the crazy idea of moving the screen to see if even larger angles ever occurred, perhaps even angles larger than 90 degrees.

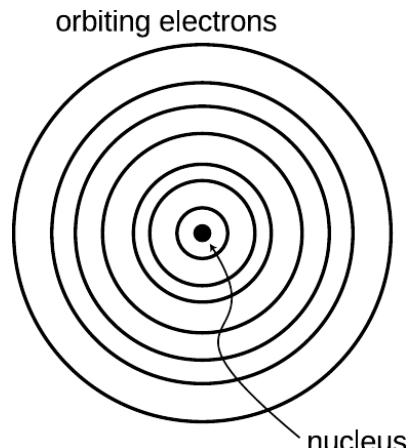
The crazy idea worked: a few alpha particles were deflected through angles of up to 180 degrees, and the routine experiment had become an epoch-making one. Rutherford said, "We have been able to get some of the alpha particles coming backwards. It was almost as incredible as if you fired a 15-inch shell at a piece of tissue paper and it came back and hit you." Explanations were hard to come by in the raisin cookie model. What intense electrical forces



g / Alpha particles being scattered by a gold nucleus. On this scale, the gold atom is the size of a car, so all the alpha particles shown here are ones that just happened to come unusually close to the nucleus. For these exceptional alpha particles, the forces from the electrons are unimportant, because they are so much more distant than the nucleus.

could have caused some of the alpha particles, moving at such astronomical speeds, to change direction so drastically? Since each gold atom was electrically neutral, it would not exert much force on an alpha particle outside it. True, if the alpha particle was very near to or inside of a particular atom, then the forces would not necessarily cancel out perfectly; if the alpha particle happened to come very close to a particular electron, the $1/r^2$ form of the Coulomb force law would make for a very strong force. But Marsden and Rutherford knew that an alpha particle was 8000 times more massive than an electron, and it is simply not possible for a more massive object to rebound backwards from a collision with a less massive object while conserving momentum and energy. It might be possible in principle for a particular alpha to follow a path that took it very close to one electron, and then very close to another electron, and so on, with the net result of a large deflection, but careful calculations showed that such multiple “close encounters” with electrons would be millions of times too rare to explain what was actually observed.

At this point, Rutherford and Marsden dusted off an unpopular and neglected model of the atom, in which all the electrons orbited around a small, positively charged core or “nucleus,” just like the planets orbiting around the sun. All the positive charge and nearly all the mass of the atom would be concentrated in the nucleus, rather than spread throughout the atom as in the raisin cookie model. The positively charged alpha particles would be repelled by the gold atom’s nucleus, but most of the alphas would not come close enough to any nucleus to have their paths drastically altered. The few that did come close to a nucleus, however, could rebound backwards from a single such encounter, since the nucleus of a heavy gold atom would be fifty times more massive than an alpha



h / The planetary model of the atom.

particle. It turned out that it was not even too difficult to derive a formula giving the relative frequency of deflections through various angles, and this calculation agreed with the data well enough (to within 15%), considering the difficulty in getting good experimental statistics on the rare, very large angles.

What had started out as a tedious exercise to get a student started in science had ended as a revolution in our understanding of nature. Indeed, the whole thing may sound a little too much like a moralistic fable of the scientific method with overtones of the Horatio Alger genre. The skeptical reader may wonder why the planetary model was ignored so thoroughly until Marsden and Rutherford's discovery. Is science really more of a sociological enterprise, in which certain ideas become accepted by the establishment, and other, equally plausible explanations are arbitrarily discarded? Some social scientists are currently ruffling a lot of scientists' feathers with critiques very much like this, but in this particular case, there were very sound reasons for rejecting the planetary model. As you'll learn in more detail later in this course, any charged particle that undergoes an acceleration dissipate energy in the form of light. In the planetary model, the electrons were orbiting the nucleus in circles or ellipses, which meant they were undergoing acceleration, just like the acceleration you feel in a car going around a curve. They should have dissipated energy as light, and eventually they should have lost all their energy. Atoms don't spontaneously collapse like that, which was why the raisin cookie model, with its stationary electrons, was originally preferred. There were other problems as well. In the planetary model, the one-electron atom would have to be flat, which would be inconsistent with the success of molecular modeling with spherical balls representing hydrogen and atoms. These molecular models also seemed to work best if specific sizes were used for different atoms, but there is no obvious reason in the planetary model why the radius of an electron's orbit should be a fixed number. In view of the conclusive Marsden-Rutherford results, however, these became fresh puzzles in atomic physics, not reasons for disbelieving the planetary model.

Some phenomena explained with the planetary model

The planetary model may not be the ultimate, perfect model of the atom, but don't underestimate its power. It already allows us to visualize correctly a great many phenomena.

As an example, let's consider the distinctions among nonmetals, metals that are magnetic, and metals that are nonmagnetic. As shown in figure i, a metal differs from a nonmetal because its outermost electrons are free to wander rather than owing their allegiance to a particular atom. A metal that can be magnetized is one that is willing to line up the rotations of some of its electrons so that their axes are parallel. Recall that magnetic forces are forces made

by moving charges; we have not yet discussed the mathematics and geometry of magnetic forces, but it is easy to see how random orientations of the atoms in the nonmagnetic substance would lead to cancellation of the forces.

Even if the planetary model does not immediately answer such questions as why one element would be a metal and another a non-metal, these ideas would be difficult or impossible to conceptualize in the raisin cookie model.

Discussion Question

A In reality, charges of the same type repel one another and charges of different types are attracted. Suppose the rules were the other way around, giving repulsion between opposite charges and attraction between similar ones. What would the universe be like?

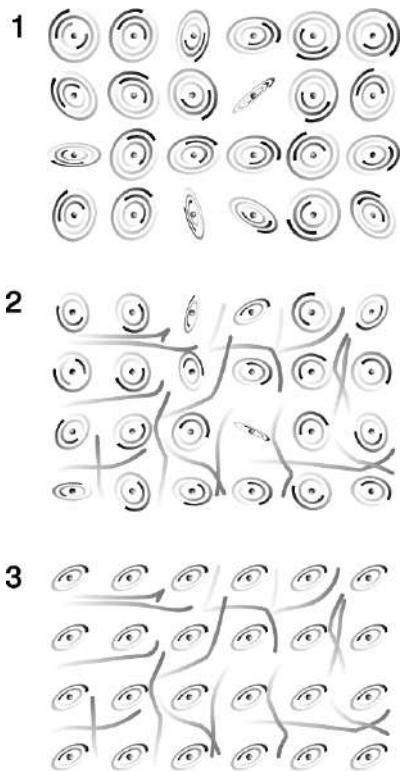
8.2.3 Atomic number

As alluded to in a discussion question in the previous section, scientists of this period had only a very approximate idea of how many units of charge resided in the nuclei of the various chemical elements. Although we now associate the number of units of nuclear charge with the element's position on the periodic table, and call it the atomic number, they had no idea that such a relationship existed. Mendeleev's table just seemed like an organizational tool, not something with any necessary physical significance. And everything Mendeleev had done seemed equally valid if you turned the table upside-down or reversed its left and right sides, so even if you wanted to number the elements sequentially with integers, there was an ambiguity as to how to do it. Mendeleev's original table was in fact upside-down compared to the modern one.

¹ H															² He		
³ Li	⁴ Be																
¹¹ Na	¹² Mg																
¹⁹ K	²⁰ Ca	²¹ Sc															
³⁷ Rb	³⁸ Sr	³⁹ Y	²² Ti	²³ V	²⁴ Cr	²⁵ Mn	²⁶ Fe	²⁷ Co	²⁸ Ni	²⁹ Cu	³⁰ Zn	³¹ Ga	³² Ge	³³ As	³⁴ Se	³⁵ Br	³⁶ Kr
		*	⁴⁰ Zr	⁴¹ Nb	⁴² Mo	⁴³ Tc	⁴⁴ Ru	⁴⁵ Rh	⁴⁶ Pd	⁴⁷ Ag	⁴⁸ Cd	⁴⁹ In	⁵⁰ Sn	⁵¹ Sb	⁵² Te	⁵³ I	⁵⁴ Xe
⁵⁵ Cs	⁵⁶ Ba	⁵⁷ La *	⁷² Hf	⁷³ Ta	⁷⁴ W	⁷⁵ Re	⁷⁶ Os	⁷⁷ Ir	⁷⁸ Pt	⁷⁹ Au	⁸⁰ Hg	⁸¹ Tl	⁸² Pb	⁸³ Bi	⁸⁴ Po	⁸⁵ At	⁸⁶ Rn
⁸⁷ Fr	⁸⁸ Ra	⁸⁹ Ac **	¹⁰⁴ Rf	¹⁰⁵ Ha	¹⁰⁶	¹⁰⁷	¹⁰⁸	¹⁰⁹	¹¹⁰	¹¹¹	¹¹²	¹¹³	¹¹⁴	¹¹⁵	¹¹⁶	¹¹⁷	¹¹⁸
*	⁵⁸ Ce	⁵⁹ Pr	⁶⁰ Nd	⁶¹ Pm	⁶² Sm	⁶³ Eu	⁶⁴ Gd	⁶⁵ Tb	⁶⁶ Dy	⁶⁷ Ho	⁶⁸ Er	⁶⁹ Tm	⁷⁰ Yb	⁷¹ Lu			
**	⁹⁰ Th	⁹¹ Pa	⁹² U	⁹³ Np	⁹⁴ Pu	⁹⁵ Am	⁹⁶ Cm	⁹⁷ Bk	⁹⁸ Cf	⁹⁹ Es	¹⁰⁰ Fm	¹⁰¹ Md	¹⁰² No	¹⁰³ Lr			

j / A modern periodic table, labeled with atomic numbers. Mendeleev's original table was upside-down compared to this one.

In the period immediately following the discovery of the nucleus, physicists only had rough estimates of the charges of the various nu-

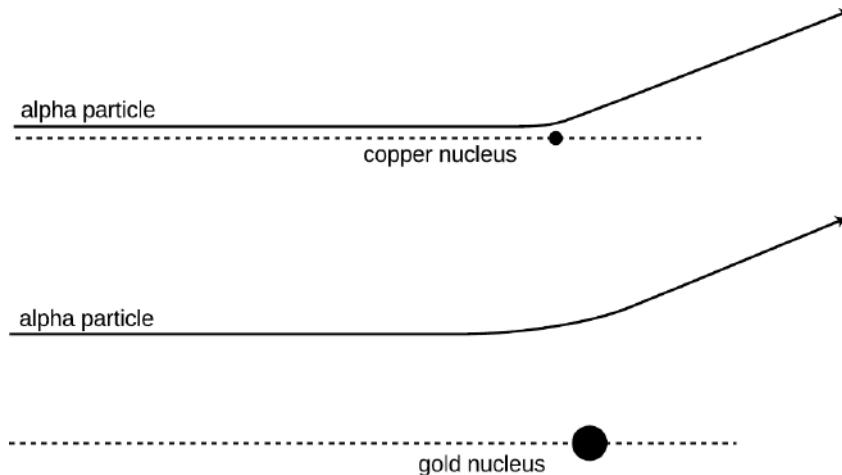


i / The planetary model applied to a nonmetal, 1, an unmagnetized metal, 2, and a magnetized metal, 3. Note that these figures are all simplified in several ways. For one thing, the electrons of an individual atom do not all revolve around the nucleus in the same plane. It is also very unusual for a metal to become so strongly magnetized that 100% of its atoms have their rotations aligned as shown in this figure.

clei. In the case of the very lightest nuclei, they simply found the maximum number of electrons they could strip off by various methods: chemical reactions, electric sparks, ultraviolet light, and so on. For example they could easily strip off one or two electrons from helium, making He^+ or He^{++} , but nobody could make He^{+++} , presumably because the nuclear charge of helium was only $+2e$. Unfortunately only a few of the lightest elements could be stripped completely, because the more electrons were stripped off, the greater the positive net charge remaining, and the more strongly the rest of the negatively charged electrons would be held on. The heavy elements' atomic numbers could only be roughly extrapolated from the light elements, where the atomic number was about half the atom's mass expressed in units of the mass of a hydrogen atom. Gold, for example, had a mass about 197 times that of hydrogen, so its atomic number was estimated to be about half that, or somewhere around 100. We now know it to be 79.

How did we finally find out? The riddle of the nuclear charges was at last successfully attacked using two different techniques, which gave consistent results. One set of experiments, involving x-rays, was performed by the young Henry Mosely, whose scientific brilliance was soon to be sacrificed in a battle between European imperialists over who would own the Dardanelles, during that pointless conflict then known as the War to End All Wars, and now referred to as World War I.

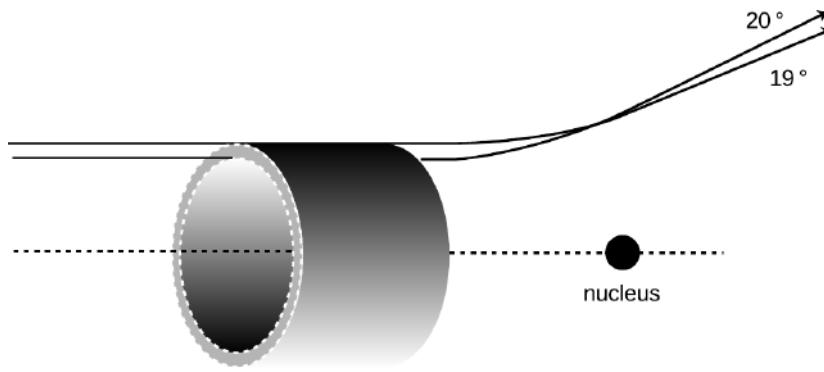
k / An alpha particle has to come much closer to the low-charged copper nucleus in order to be deflected through the same angle.



Since Mosely's analysis requires several concepts with which you are not yet familiar, we will instead describe the technique used by James Chadwick at around the same time. An added bonus of describing Chadwick's experiments is that they presaged the important modern technique of studying *collisions* of subatomic particles. In grad school, I worked with a professor whose thesis adviser's thesis adviser was Chadwick, and he related some interesting stories about the man. Chadwick was apparently a little nutty and a com-

plete fanatic about science, to the extent that when he was held in a German prison camp during World War II, he managed to cajole his captors into allowing him to scrounge up parts from broken radios so that he could attempt to do physics experiments.

Chadwick's experiment worked like this. Suppose you perform two Rutherford-type alpha scattering measurements, first one with a gold foil as a target as in Rutherford's original experiment, and then one with a copper foil. It is possible to get large angles of deflection in both cases, but as shown in figure 1, the alpha particle must be heading almost straight for the copper nucleus to get the same angle of deflection that would have occurred with an alpha that was much farther off the mark; the gold nucleus' charge is so much greater than the copper's that it exerts a strong force on the alpha particle even from far off. The situation is very much like that of a blindfolded person playing darts. Just as it is impossible to aim an alpha particle at an individual nucleus in the target, the blindfolded person cannot really aim the darts. Achieving a very close encounter with the copper atom would be akin to hitting an inner circle on the dartboard. It's much more likely that one would have the luck to hit the outer circle, which covers a greater number of square inches. By analogy, if you measure the frequency with which alphas are scattered by copper at some particular angle, say between 19 and 20 degrees, and then perform the same measurement at the same angle with gold, you get a much higher percentage for gold than for copper.



I / An alpha particle must be headed for the ring on the front of the imaginary cylindrical pipe in order to produce scattering at an angle between 19 and 20 degrees. The area of this ring is called the "cross-section" for scattering at 19-20° because it is the cross-sectional area of a cut through the pipe.

In fact, the numerical ratio of the two nuclei's charges can be derived from this same experimentally determined ratio. Using the standard notation Z for the atomic number (charge of the nucleus divided by e), the following equation can be proved (example 6):

$$\frac{Z_{gold}^2}{Z_{copper}^2} = \frac{\text{number of alphas scattered by gold at } 19-20^\circ}{\text{number of alphas scattered by copper at } 19-20^\circ}$$

By making such measurements for targets constructed from all the elements, one can infer the ratios of all the atomic numbers, and

since the atomic numbers of the light elements were already known, atomic numbers could be assigned to the entire periodic table. According to Moseley, the atomic numbers of copper, silver and platinum were 29, 47, and 78, which corresponded well with their positions on the periodic table. Chadwick's figures for the same elements were 29.3, 46.3, and 77.4, with error bars of about 1.5 times the fundamental charge, so the two experiments were in good agreement.

The point here is absolutely not that you should be ready to plug numbers into the above equation for a homework or exam question! My overall goal in this chapter is to explain how we know what we know about atoms. An added bonus of describing Chadwick's experiment is that the approach is very similar to that used in modern particle physics experiments, and the ideas used in the analysis are closely related to the now-ubiquitous concept of a "cross-section." In the dartboard analogy, the cross-section would be the area of the circular ring you have to hit. The reasoning behind the invention of the term "cross-section" can be visualized as shown in figure 1. In this language, Rutherford's invention of the planetary model came from his unexpected discovery that there was a nonzero cross-section for alpha scattering from gold at large angles, and Chadwick confirmed Moseley's determinations of the atomic numbers by measuring cross-sections for alpha scattering.

Proof of the relationship between Z and scattering example 6

The equation above can be derived by the following not very rigorous proof. To deflect the alpha particle by a certain angle requires that it acquire a certain momentum component in the direction perpendicular to its original momentum. Although the nucleus's force on the alpha particle is not constant, we can pretend that it is approximately constant during the time when the alpha is within a distance equal to, say, 150% of its distance of closest approach, and that the force is zero before and after that part of the motion. (If we chose 120% or 200%, it shouldn't make any difference in the final result, because the final result is a ratio, and the effects on the numerator and denominator should cancel each other.) In the approximation of constant force, the change in the alpha's perpendicular momentum component is then equal to $F\Delta t$. The Coulomb force law says the force is proportional to Z/r^2 . Although r does change somewhat during the time interval of interest, it's good enough to treat it as a constant number, since we're only computing the ratio between the two experiments' results. Since we are approximating the force as acting over the time during which the distance is not too much greater than the distance of closest approach, the time interval Δt must be proportional to r , and the sideways momentum imparted to the alpha, $F\Delta t$, is proportional to $(Z/r^2)r$, or Z/r . If we're comparing alphas

scattered at the same angle from gold and from copper, then Δp is the same in both cases, and the proportionality $\Delta p \propto Z/r$ tells us that the ones scattered from copper at that angle had to be headed in along a line closer to the central axis by a factor equaling $Z_{\text{gold}}/Z_{\text{copper}}$. If you imagine a “dartboard ring” that the alphas have to hit, then the ring for the gold experiment has the same proportions as the one for copper, but it is enlarged by a factor equal to $Z_{\text{gold}}/Z_{\text{copper}}$. That is, not only is the radius of the ring greater by that factor, but unlike the rings on a normal dartboard, the thickness of the outer ring is also greater in proportion to its radius. When you take a geometric shape and scale it up in size like a photographic enlargement, its area is increased in proportion to the square of the enlargement factor, so the area of the dartboard ring in the gold experiment is greater by a factor equal to $(Z_{\text{gold}}/Z_{\text{copper}})^2$. Since the alphas are aimed entirely randomly, the chances of an alpha hitting the ring are in proportion to the area of the ring, which proves the equation given above.

As an example of the modern use of scattering experiments and cross-section measurements, you may have heard of the recent experimental evidence for the existence of a particle called the top quark. Of the twelve subatomic particles currently believed to be the smallest constituents of matter, six form a family called the quarks, distinguished from the other six by the intense attractive forces that make the quarks stick to each other. (The other six consist of the electron plus five other, more exotic particles.) The only two types of quarks found in naturally occurring matter are the “up quark” and “down quark,” which are what protons and neutrons are made of, but four other types were theoretically predicted to exist, for a total of six. (The whimsical term “quark” comes from a line by James Joyce reading “Three quarks for master Mark.”) Until recently, only five types of quarks had been proven to exist via experiments, and the sixth, the top quark, was only theorized. There was no hope of ever detecting a top quark directly, since it is radioactive, and only exists for a zillionth of a second before evaporating. Instead, the researchers searching for it at the Fermi National Accelerator Laboratory near Chicago measured cross-sections for scattering of nuclei off of other nuclei. The experiment was much like those of Rutherford and Chadwick, except that the incoming nuclei had to be boosted to much higher speeds in a particle accelerator. The resulting encounter with a target nucleus was so violent that both nuclei were completely demolished, but, as Einstein proved, energy can be converted into matter, and the energy of the collision creates a spray of exotic, radioactive particles, like the deadly shower of wood fragments produced by a cannon ball in an old naval battle. Among those particles were some top quarks. The cross-sections being measured were the cross-sections for the production of certain combinations of these secondary particles. However different the

details, the principle was the same as that employed at the turn of the century: you smash things together and look at the fragments that fly off to see what was inside them. The approach has been compared to shooting a clock with a rifle and then studying the pieces that fly off to figure out how the clock worked.

Discussion Questions

A The diagram, showing alpha particles being deflected by a gold nucleus, was drawn with the assumption that alpha particles came in on lines at many different distances from the nucleus. Why wouldn't they all come in along the same line, since they all came out through the same tube?

B Why does it make sense that, as shown in the figure, the trajectories that result in 19° and 20° scattering cross each other?

C Rutherford knew the velocity of the alpha particles emitted by radium, and guessed that the positively charged part of a gold atom had a charge of about $+100e$ (we now know it is $+79e$). Considering the fact that some alpha particles were deflected by 180° , how could he then use conservation of energy to derive an upper limit on the size of a gold nucleus? (For simplicity, assume the size of the alpha particle is negligible compared to that of the gold nucleus, and ignore the fact that the gold nucleus recoils a little from the collision, picking up a little kinetic energy.)

8.2.4 The structure of nuclei

The proton

The fact that the nuclear charges were all integer multiples of e suggested to many physicists that rather than being a pointlike object, the nucleus might contain smaller particles having individual charges of $+e$. Evidence in favor of this idea was not long in arriving. Rutherford reasoned that if he bombarded the atoms of a very light element with alpha particles, the small charge of the target nuclei would give a very weak repulsion. Perhaps those few alpha particles that happened to arrive on head-on collision courses would get so close that they would physically crash into some of the target nuclei. An alpha particle is itself a nucleus, so this would be a collision between two nuclei, and a violent one due to the high speeds involved. Rutherford hit pay dirt in an experiment with alpha particles striking a target containing nitrogen atoms. Charged particles were detected flying out of the target like parts flying off of cars in a high-speed crash. Measurements of the deflection of these particles in electric and magnetic fields showed that they had the same charge-to-mass ratio as singly-ionized hydrogen atoms. Rutherford concluded that these were the conjectured singly-charged particles that held the charge of the nucleus, and they were later named protons. The hydrogen nucleus consists of a single proton, and in general, an element's atomic number gives the number of protons contained in each of its nuclei. The mass of the proton is about 1800 times greater than the mass of the electron.

The neutron

It would have been nice and simple if all the nuclei could have been built only from protons, but that couldn't be the case. If you spend a little time looking at a periodic table, you will soon notice that although some of the atomic masses are very nearly integer multiples of hydrogen's mass, many others are not. Even where the masses are close whole numbers, the masses of an element other than hydrogen is always greater than its atomic number, not equal to it. Helium, for instance, has two protons, but its mass is four times greater than that of hydrogen.

Chadwick cleared up the confusion by proving the existence of a new subatomic particle. Unlike the electron and proton, which are electrically charged, this particle is electrically neutral, and he named it the neutron. Chadwick's experiment has been described in detail on p. 140, but briefly the method was to expose a sample of the light element beryllium to a stream of alpha particles from a lump of radium. Beryllium has only four protons, so an alpha that happens to be aimed directly at a beryllium nucleus can actually hit it rather than being stopped short of a collision by electrical repulsion. Neutrons were observed as a new form of radiation emerging from the collisions, and Chadwick correctly inferred that they were previously unsuspected components of the nucleus that had been knocked out. As described earlier, Chadwick also determined the mass of the neutron; it is very nearly the same as that of the proton.

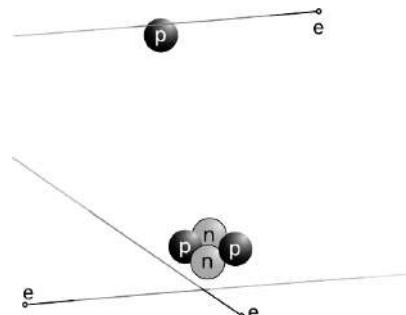
To summarize, atoms are made of three types of particles:

	charge	mass in units of the proton's mass	location in atom
proton	$+e$	1	in nucleus
neutron	0	1.001	in nucleus
electron	$-e$	1/1836	orbiting nucleus

The existence of neutrons explained the mysterious masses of the elements. Helium, for instance, has a mass very close to four times greater than that of hydrogen. This is because it contains two neutrons in addition to its two protons. The mass of an atom is essentially determined by the total number of neutrons and protons. The total number of neutrons plus protons is therefore referred to as the atom's *mass number*.

Isotopes

We now have a clear interpretation of the fact that helium is close to four times more massive than hydrogen, and similarly for all the atomic masses that are close to an integer multiple of the mass of hydrogen. But what about copper, for instance, which had an atomic mass 63.5 times that of hydrogen? It didn't seem reasonable to think that it possessed an extra half of a neutron! The



m / Examples of the construction of atoms: hydrogen (top) and helium (bottom). On this scale, the electrons' orbits would be the size of a college campus.

solution was found by measuring the mass-to-charge ratios of singly-ionized atoms (atoms with one electron removed). The technique is essentially that same as the one used by Thomson for cathode rays, except that whole atoms do not spontaneously leap out of the surface of an object as electrons sometimes do. Figure n shows an example of how the ions can be created and injected between the charged plates for acceleration.

Injecting a stream of copper ions into the device, we find a surprise — the beam splits into two parts! Chemists had elevated to dogma the assumption that all the atoms of a given element were identical, but we find that 69% of copper atoms have one mass, and 31% have another. Not only that, but both masses are very nearly integer multiples of the mass of hydrogen (63 and 65, respectively). Copper gets its chemical identity from the number of protons in its nucleus, 29, since chemical reactions work by electric forces. But apparently some copper atoms have $63 - 29 = 34$ neutrons while others have $65 - 29 = 36$. The atomic mass of copper, 63.5, reflects the proportions of the mixture of the mass-63 and mass-65 varieties. The different mass varieties of a given element are called *isotopes* of that element.

Isotopes can be named by giving the mass number as a subscript to the left of the chemical symbol, e.g., ^{65}Cu . Examples:

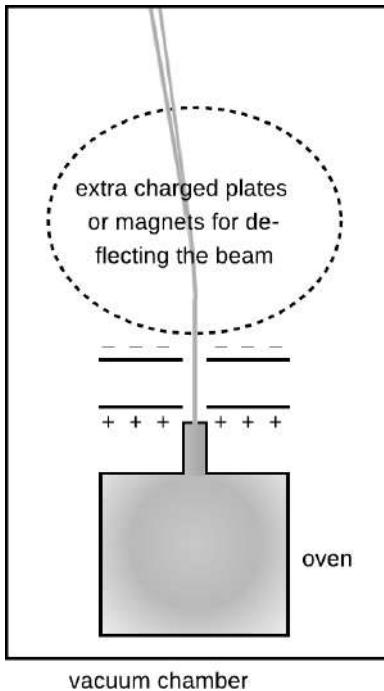
	<i>protons</i>	<i>neutrons</i>	<i>mass number</i>
^1H	1	0	$0+1 = 1$
^4He	2	2	$2+2 = 4$
^{12}C	6	6	$6+6 = 12$
^{14}C	6	8	$6+8 = 14$
^{262}Ha	105	157	$105+157 = 262$

self-check D

Why are the positive and negative charges of the accelerating plates reversed in the isotope-separating apparatus compared to the Thomson apparatus?

▷ Answer, p. 1062

Chemical reactions are all about the exchange and sharing of electrons: the nuclei have to sit out this dance because the forces of electrical repulsion prevent them from ever getting close enough to make contact with each other. Although the protons do have a vitally important effect on chemical processes because of their electrical forces, the neutrons can have no effect on the atom's chemical reactions. It is not possible, for instance, to separate ^{63}Cu from ^{65}Cu by chemical reactions. This is why chemists had never realized that different isotopes existed. (To be perfectly accurate, different isotopes do behave slightly differently because the more massive atoms move more sluggishly and therefore react with a tiny bit less intensity. This tiny difference is used, for instance, to separate out the isotopes of uranium needed to build a nuclear bomb. The smallness of this effect makes the separation process a slow and difficult one,



n / A version of the Thomson apparatus modified for measuring the mass-to-charge ratios of ions rather than electrons. A small sample of the element in question, copper in our example, is boiled in the oven to create a thin vapor. (A vacuum pump is continuously sucking on the main chamber to keep it from accumulating enough gas to stop the beam of ions.) Some of the atoms of the vapor are ionized by a spark or by ultraviolet light. Ions that wander out of the nozzle and into the region between the charged plates are then accelerated toward the top of the figure. As in the Thomson experiment, mass-to-charge ratios are inferred from the deflection of the beam.

which is what we have to thank for the fact that nuclear weapons have not been built by every terrorist cabal on the planet. See also example 1, p. 965.)

Sizes and shapes of nuclei

Matter is nearly all nuclei if you count by weight, but in terms of volume nuclei don't amount to much. The radius of an individual neutron or proton is very close to 1 fm ($1 \text{ fm} = 10^{-15} \text{ m}$), so even a big lead nucleus with a mass number of 208 still has a diameter of only about 13 fm, which is ten thousand times smaller than the diameter of a typical atom. Contrary to the usual imagery of the nucleus as a small sphere, it turns out that many nuclei are somewhat elongated, like an American football, and a few have exotic asymmetric shapes like pears or kiwi fruits.

Discussion Questions

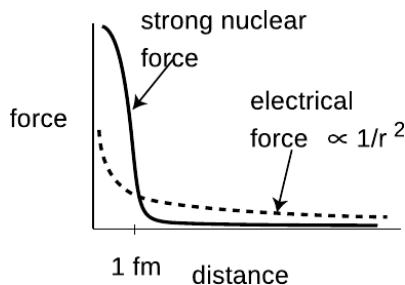
A Suppose the entire universe was in a (very large) cereal box, and the nutritional labeling was supposed to tell a godlike consumer what percentage of the contents was nuclei. Roughly what would the percentage be like if the labeling was according to mass? What if it was by volume?



o / A nuclear power plant at Cattenom, France. Unlike the coal and oil plants that supply most of the U.S.'s electrical power, a nuclear power plant like this one releases no pollution or greenhouse gases into the Earth's atmosphere, and therefore doesn't contribute to global warming. The white stuff puffing out of this plant is non-radioactive water vapor. Although nuclear power plants generate long-lived nuclear waste, this waste arguably poses much less of a threat to the biosphere than greenhouse gases would.

8.2.5 The strong nuclear force, alpha decay and fission

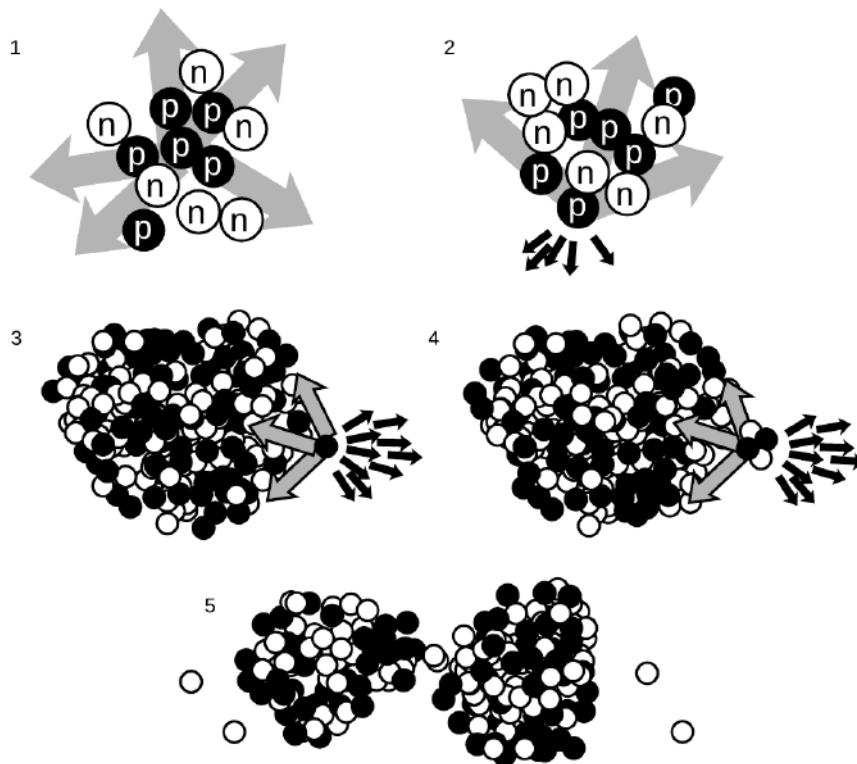
Once physicists realized that nuclei consisted of positively charged protons and uncharged neutrons, they had a problem on their hands. The electrical forces among the protons are all repulsive, so the nucleus should simply fly apart! The reason all the nuclei in your body are not spontaneously exploding at this moment is that there is another force acting. This force, called the *strong nuclear force*, is



p / The strong nuclear force cuts off very sharply at a range of about 1 fm.

always attractive, and acts between neutrons and neutrons, neutrons and protons, and protons and protons with roughly equal strength. The strong nuclear force does not have any effect on electrons, which is why it does not influence chemical reactions.

Unlike electric forces, whose strengths are given by the simple Coulomb force law, there is no simple formula for how the strong nuclear force depends on distance. Roughly speaking, it is effective over ranges of ~ 1 fm, but falls off extremely quickly at larger distances (much faster than $1/r^2$). Since the radius of a neutron or proton is about 1 fm, that means that when a bunch of neutrons and protons are packed together to form a nucleus, the strong nuclear force is effective only between neighbors.



q / 1. The forces cancel. 2. The forces don't cancel. 3. In a heavy nucleus, the large number of electrical repulsions can add up to a force that is comparable to the strong nuclear attraction. 4. Alpha emission. 5. Fission.

Figure q illustrates how the strong nuclear force acts to keep ordinary nuclei together, but is not able to keep very heavy nuclei from breaking apart. In q/1, a proton in the middle of a carbon nucleus feels an attractive strong nuclear force (arrows) from each of its nearest neighbors. The forces are all in different directions, and tend to cancel out. The same is true for the repulsive electrical forces (not shown). In figure q/2, a proton at the edge of the nucleus has neighbors only on one side, and therefore all the strong nuclear

forces acting on it are tending to pull it back in. Although all the electrical forces from the other five protons (dark arrows) are all pushing it out of the nucleus, they are not sufficient to overcome the strong nuclear forces.

In a very heavy nucleus, $q/3$, a proton that finds itself near the edge has only a few neighbors close enough to attract it significantly via the strong nuclear force, but every other proton in the nucleus exerts a repulsive electrical force on it. If the nucleus is large enough, the total electrical repulsion may be sufficient to overcome the attraction of the strong force, and the nucleus may spit out a proton. Proton emission is fairly rare, however; a more common type of radioactive decay¹ in heavy nuclei is alpha decay, shown in $q/4$. The imbalance of the forces is similar, but the chunk that is ejected is an alpha particle (two protons and two neutrons) rather than a single proton.

It is also possible for the nucleus to split into two pieces of roughly equal size, $q/5$, a process known as fission. Note that in addition to the two large fragments, there is a spray of individual neutrons. In a nuclear fission bomb or a nuclear fission reactor, some of these neutrons fly off and hit other nuclei, causing them to undergo fission as well. The result is a chain reaction.

When a nucleus is able to undergo one of these processes, it is said to be radioactive, and to undergo radioactive decay. Some of the naturally occurring nuclei on earth are radioactive. The term “radioactive” comes from Becquerel’s image of rays radiating out from something, not from radio waves, which are a whole different phenomenon. The term “decay” can also be a little misleading, since it implies that the nucleus turns to dust or simply disappears – actually it is splitting into two new nuclei with the same total number of neutrons and protons, so the term “radioactive transformation” would have been more appropriate. Although the original atom’s electrons are mere spectators in the process of weak radioactive decay, we often speak loosely of “radioactive atoms” rather than “radioactive nuclei.”

Randomness in physics

How does an atom decide when to decay? We might imagine that it is like a termite-infested house that gets weaker and weaker, until finally it reaches the day on which it is destined to fall apart. Experiments, however, have not succeeded in detecting such “ticking clock” hidden below the surface; the evidence is that all atoms of a given isotope are absolutely identical. Why, then, would one uranium atom decay today while another lives for another million years? The answer appears to be that it is entirely random. We

¹Alpha decay is more common because an alpha particle happens to be a very stable arrangement of protons and neutrons.

can make general statements about the average time required for a certain isotope to decay, or how long it will take for half the atoms in a sample to decay (its half-life), but we can never predict the behavior of a particular atom.

This is the first example we have encountered of an inescapable randomness in the laws of physics. If this kind of randomness makes you uneasy, you're in good company. Einstein's famous quote is "...I am convinced that He [God] does not play dice." Einstein's distaste for randomness, and his association of determinism with divinity, goes back to the Enlightenment conception of the universe as a gigantic piece of clockwork that only had to be set in motion initially by the Builder. Physics had to be entirely rebuilt in the 20th century to incorporate the fundamental randomness of physics, and this modern revolution is the topic of

chapter 13. In

particular, we will delay the mathematical development of the half-life concept until then.

8.2.6 The weak nuclear force; beta decay

All the nuclear processes we've discussed so far have involved rearrangements of neutrons and protons, with no change in the total number of neutrons or the total number of protons. Now consider the proportions of neutrons and protons in your body and in the planet earth: neutrons and protons are roughly equally numerous in your body's carbon and oxygen nuclei, and also in the nickel and iron that make up most of the earth. The proportions are about 50-50. But, as discussed in more detail on p. 523, the only chemical elements produced in any significant quantities by the big bang² were hydrogen (about 90%) and helium (about 10%). If the early universe was almost nothing but hydrogen atoms, whose nuclei are protons, where did all those neutrons come from?

The answer is that there is another nuclear force, the weak nuclear force, that is capable of transforming neutrons into protons and vice-versa. Two possible reactions are



and



(There is also a third type called electron capture, in which a proton grabs one of the atom's electrons and they produce a neutron and a neutrino.)

Whereas alpha decay and fission are just a redivision of the previously existing particles, these reactions involve the destruction of

²The evidence for the big bang theory of the origin of the universe was discussed on p. 370.

one particle and the creation of three new particles that did not exist before.

There are three new particles here that you have never previously encountered. The symbol e^+ stands for an antielectron, which is a particle just like the electron in every way, except that its electric charge is positive rather than negative. Antielectrons are also known as positrons. Nobody knows why electrons are so common in the universe and antielectrons are scarce. When an antielectron encounters an electron, they annihilate each other, producing gamma rays, and this is the fate of all the antielectrons that are produced by natural radioactivity on earth. Antielectrons are an example of antimatter. A complete atom of antimatter would consist of antiprotons, antielectrons, and antineutrons. Although individual particles of antimatter occur commonly in nature due to natural radioactivity and cosmic rays, only a few complete atoms of antihydrogen have ever been produced artificially.

The notation ν stands for a particle called a neutrino, and $\bar{\nu}$ means an antineutrino. Neutrinos and antineutrinos have no electric charge (hence the name).

We can now list all four of the known fundamental forces of physics:

- gravity
- electromagnetism
- strong nuclear force
- weak nuclear force

The other forces we have learned about, such as friction and the normal force, all arise from electromagnetic interactions between atoms, and therefore are not considered to be fundamental forces of physics.

Decay of ^{212}Pb

example 7

As an example, consider the radioactive isotope of lead ^{212}Pb . It contains 82 protons and 130 neutrons. It decays by the process $n \rightarrow p + e^- + \bar{\nu}$. The newly created proton is held inside the nucleus by the strong nuclear force, so the new nucleus contains 83 protons and 129 neutrons. Having 83 protons makes it the element bismuth, so it will be an atom of ^{212}Bi .

In a reaction like this one, the electron flies off at high speed (typically close to the speed of light), and the escaping electrons are the things that make large amounts of this type of radioactivity dangerous. The outgoing electron was the first thing that tipped off scientists in the early 1900s to the existence of this type of radioactivity. Since they didn't know that the outgoing particles were

electrons, they called them beta particles, and this type of radioactive decay was therefore known as beta decay. A clearer but less common terminology is to call the two processes electron decay and positron decay.

The neutrino or antineutrino emitted in such a reaction pretty much ignores all matter, because its lack of charge makes it immune to electrical forces, and it also remains aloof from strong nuclear interactions. Even if it happens to fly off going straight down, it is almost certain to make it through the entire earth without interacting with any atoms in any way. It ends up flying through outer space forever. The neutrino's behavior makes it exceedingly difficult to detect, and when beta decay was first discovered nobody realized that neutrinos even existed. We now know that the neutrino carries off some of the energy produced in the reaction, but at the time it seemed that the total energy afterwards (not counting the unsuspected neutrino's energy) was greater than the total energy before the reaction, violating conservation of energy. Physicists were getting ready to throw conservation of energy out the window as a basic law of physics when indirect evidence led them to the conclusion that neutrinos existed.

Discussion Questions

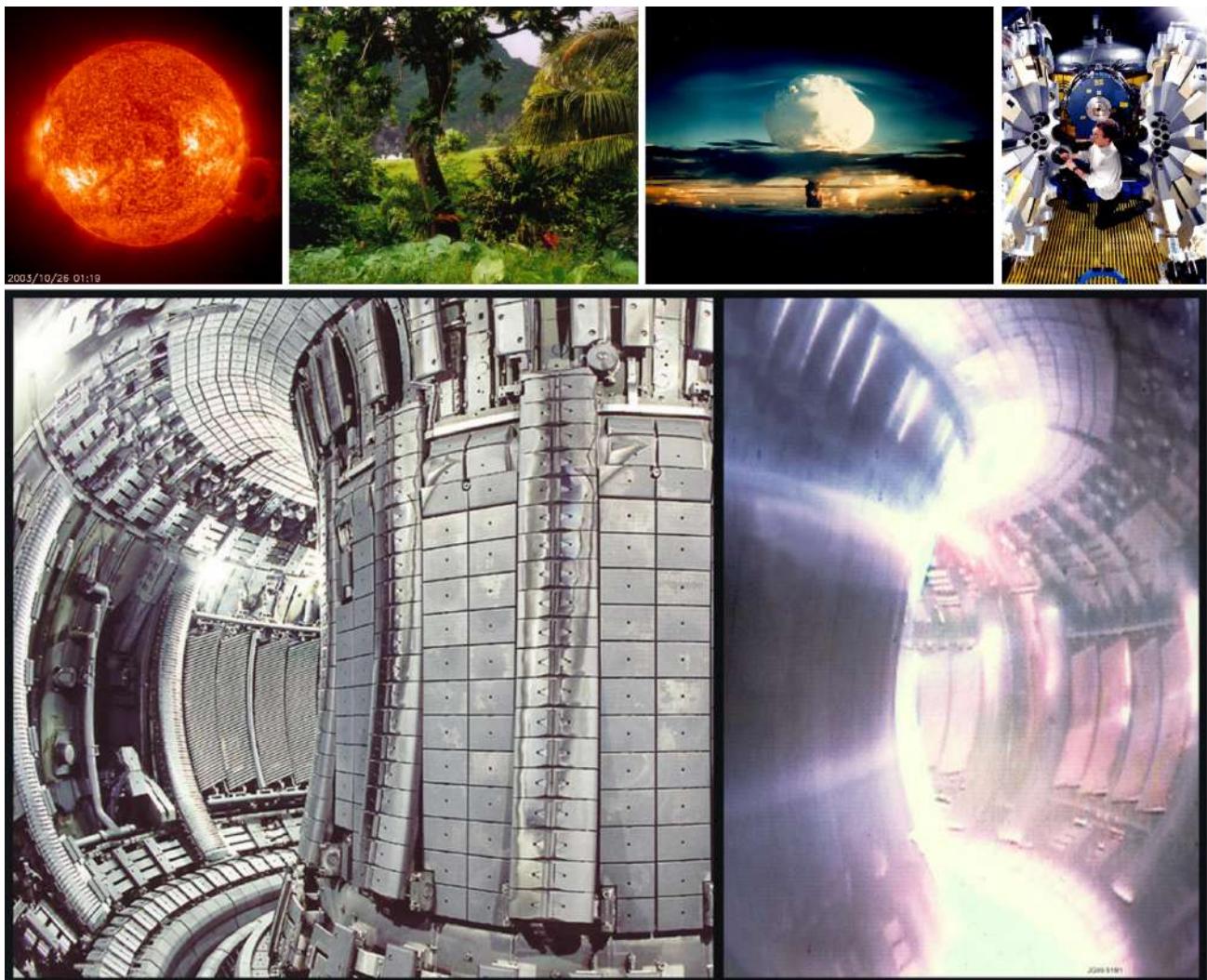
A In the reactions $n \rightarrow p + e^- + \bar{\nu}$ and $p \rightarrow n + e^+ + \nu$, verify that charge is conserved. In beta decay, when one of these reactions happens to a neutron or proton within a nucleus, one or more gamma rays may also be emitted. Does this affect conservation of charge? Would it be possible for some extra electrons to be released without violating charge conservation?

B When an antielectron and an electron annihilate each other, they produce two gamma rays. Is charge conserved in this reaction?

8.2.7 Fusion

As we have seen, heavy nuclei tend to fly apart because each proton is being repelled by every other proton in the nucleus, but is only attracted by its nearest neighbors. The nucleus splits up into two parts, and as soon as those two parts are more than about 1 fm apart, the strong nuclear force no longer causes the two fragments to attract each other. The electrical repulsion then accelerates them, causing them to gain a large amount of kinetic energy. This release of kinetic energy is what powers nuclear reactors and fission bombs.

It might seem, then, that the lightest nuclei would be the most stable, but that is not the case. Let's compare an extremely light nucleus like ${}^4\text{He}$ with a somewhat heavier one, ${}^{16}\text{O}$. A neutron or proton in ${}^4\text{He}$ can be attracted by the three others, but in ${}^{16}\text{O}$, it might have five or six neighbors attracting it. The ${}^{16}\text{O}$ nucleus is therefore more stable.



r / 1. Our sun's source of energy is nuclear fusion, so nuclear fusion is also the source of power for all life on earth, including, 2, this rain forest in Fatu-Hiva. 3. The first release of energy by nuclear fusion through human technology was the 1952 Ivy Mike test at the Enewetak Atoll. 4. This array of gamma-ray detectors is called GAMMASPHERE. During operation, the array is closed up, and a beam of ions produced by a particle accelerator strikes a target at its center, producing nuclear fusion reactions. The gamma rays can be studied for information about the structure of the fused nuclei, which are typically varieties not found in nature. 5. Nuclear fusion promises to be a clean, inexhaustible source of energy. However, the goal of commercially viable nuclear fusion power has remained elusive, due to the engineering difficulties involved in magnetically containing a plasma (ionized gas) at a sufficiently high temperature and density. This photo shows the experimental JET reactor, with the device opened up on the left, and in action on the right.

It turns out that the most stable nuclei of all are those around nickel and iron, having about 30 protons and 30 neutrons. Just as a nucleus that is too heavy to be stable can release energy by splitting apart into pieces that are closer to the most stable size, light nuclei can release energy if you stick them together to make bigger nuclei that are closer to the most stable size. Fusing one nucleus with another is called nuclear fusion. Nuclear fusion is what powers our

sun and other stars.

8.2.8 Nuclear energy and binding energies

In the same way that chemical reactions can be classified as exothermic (releasing energy) or endothermic (requiring energy to react), so nuclear reactions may either release or use up energy. The energies involved in nuclear reactions are greater by a huge factor. Thousands of tons of coal would have to be burned to produce as much energy as would be produced in a nuclear power plant by one kg of fuel.

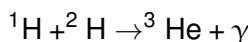
Although nuclear reactions that use up energy (endothermic reactions) can be initiated in accelerators, where one nucleus is rammed into another at high speed, they do not occur in nature, not even in the sun. The amount of kinetic energy required is simply not available.

To find the amount of energy consumed or released in a nuclear reaction, you need to know how much nuclear interaction energy, U_{nuc} , was stored or released. Experimentalists have determined the amount of nuclear energy stored in the nucleus of every stable element, as well as many unstable elements. This is the amount of mechanical work that would be required to pull the nucleus apart into its individual neutrons and protons, and is known as the nuclear binding energy.

A reaction occurring in the sun

example 8

The sun produces its energy through a series of nuclear fusion reactions. One of the reactions is



The excess energy is almost all carried off by the gamma ray (not by the kinetic energy of the helium-3 atom). The binding energies in units of pJ (picojoules) are:

$$^1\text{H} \quad 0 \text{ J}$$

$$^2\text{H} \quad 0.35593 \text{ pJ}$$

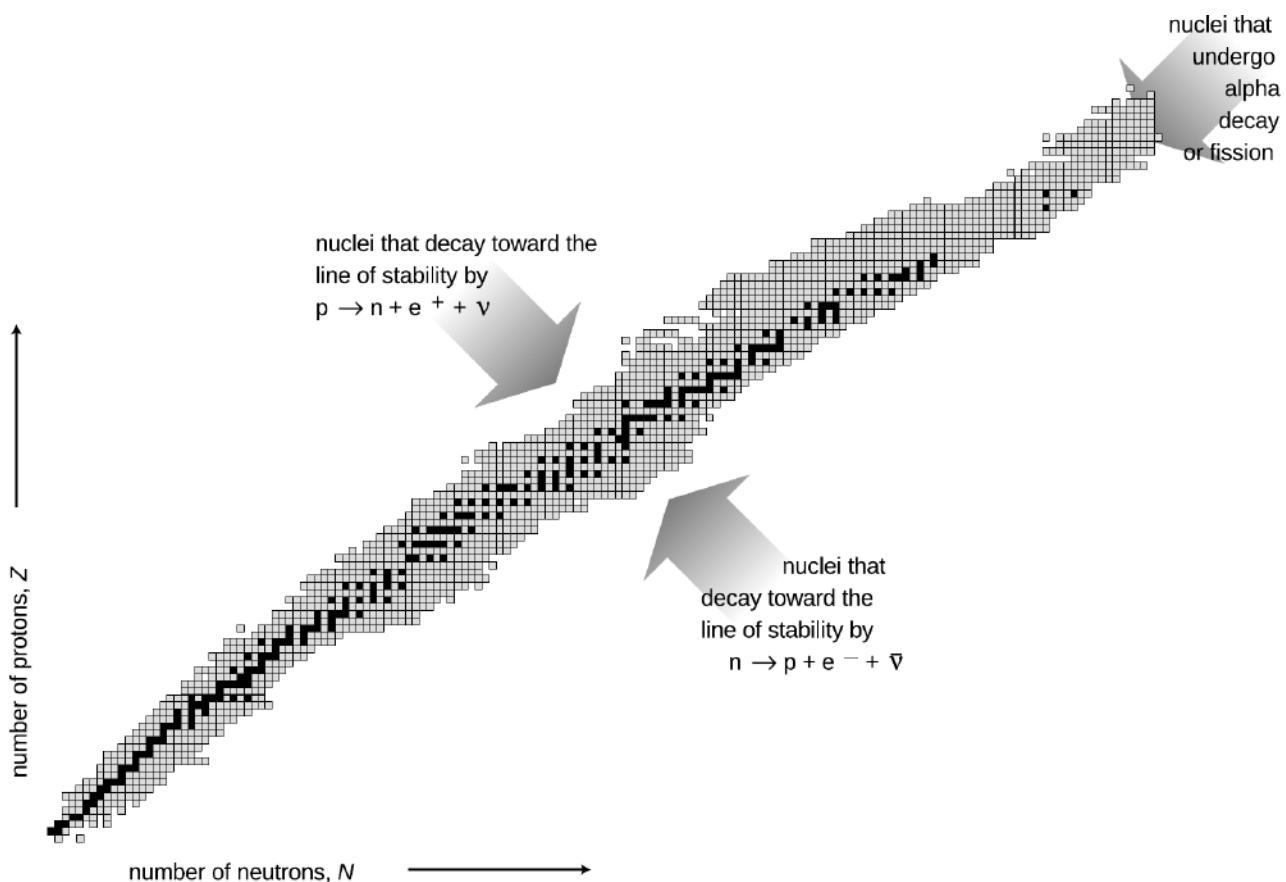
$$^3\text{He} \quad 1.23489 \text{ pJ}$$

The total initial nuclear energy is 0 pJ+0.35593 pJ, and the final nuclear energy is 1.23489 pJ, so by conservation of energy, the gamma ray must carry off 0.87896 pJ of energy. The gamma ray is then absorbed by the sun and converted to heat.

self-check E

Why is the binding energy of ^1H exactly equal to zero? ▷ Answer, p. 1062

Figure s is a compact way of showing the vast variety of the nuclei. Each box represents a particular number of neutrons and protons. The black boxes are nuclei that are stable, i.e., that would require an input of energy in order to change into another. The



s / The known nuclei, represented on a chart of proton number versus neutron number. Note the two nuclei in the bottom row with zero protons.



t / A map showing levels of radiation near the site of the Chernobyl nuclear accident.

gray boxes show all the unstable nuclei that have been studied experimentally. Some of these last for billions of years on the average before decaying and are found in nature, but most have much shorter average lifetimes, and can only be created and studied in the laboratory.

The curve along which the stable nuclei lie is called the line of stability. Nuclei along this line have the most stable proportion of neutrons to protons. For light nuclei the most stable mixture is about 50-50, but we can see that stable heavy nuclei have two or three times more neutrons than protons. This is because the electrical repulsions of all the protons in a heavy nucleus add up to a powerful force that would tend to tear it apart. The presence of a large number of neutrons increases the distances among the protons, and also increases the number of attractions due to the strong nuclear force.

8.2.9 Biological effects of ionizing radiation

Units used to measure exposure

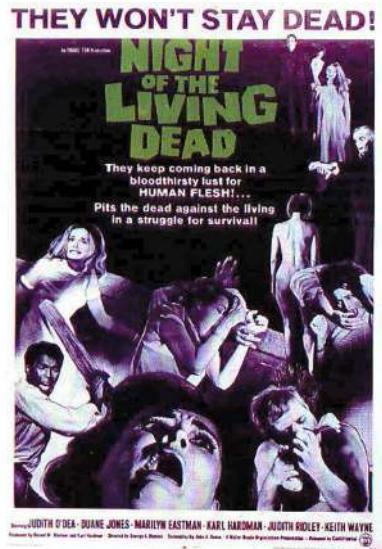
As a science educator, I find it frustrating that nowhere in the massive amount of journalism devoted to nuclear safety does one ever find any numerical statements about the amount of radiation to which people have been exposed. Anyone capable of understanding sports statistics or weather reports ought to be able to understand such measurements, as long as something like the following explanatory text was inserted somewhere in the article:

Radiation exposure is measured in units of Sieverts (Sv). The average person is exposed to about 2000 μSv (microSieverts) each year from natural background sources.

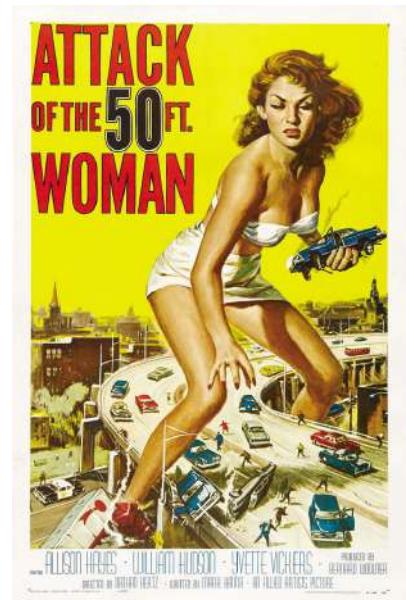
With this context, people would be able to come to informed conclusions. For example, figure t shows a scary-looking map of the levels of radiation in the area surrounding the 1986 nuclear accident at Chernobyl, Ukraine, the most serious that has ever occurred. At the boundary of the most highly contaminated (bright red) areas, people would be exposed to about 13,000 μSv per year, or about four times the natural background level. In the pink areas, which are still densely populated, the exposure is comparable to the natural level found in a high-altitude city such as Denver.

What is a Sievert? It measures the amount of energy per kilogram deposited in the body by ionizing radiation, multiplied by a “quality factor” to account for the different health hazards posed by alphas, betas, gammas, neutrons, and other types of radiation. Only ionizing radiation is counted, since nonionizing radiation simply heats one’s body rather than killing cells or altering DNA. For instance, alpha particles are typically moving so fast that their kinetic energy is sufficient to ionize thousands of atoms, but it is possible for an alpha particle to be moving so slowly that it would not have enough kinetic energy to ionize even one atom.

Unfortunately, most people don’t know much about radiation and tend to react to it based on unscientific cultural notions. These may, as in figure u, be based on fictional tropes silly enough to require the suspension of disbelief by the audience, but they can also be more subtle. People of my kids’ generation are more familiar with the 2011 Fukushima nuclear accident than with the much more serious Chernobyl accident. The news coverage of Fukushima showed scary scenes of devastated landscapes and distraught evacuees, implying that people had been killed and displaced by the release of radiation from the reaction. In fact, there were no deaths at all due to the radiation released at Fukushima, and no excess cancer deaths are statistically predicted in the future. The devastation and the death toll of 16,000 were caused by the earthquake and tsunami, which were also what damaged the plant.



u / In this classic zombie flick, a newscaster speculates that the dead have been reanimated due to radiation brought back to earth by a space probe.



v / Radiation doesn’t mutate entire multicellular organisms.

Effects of exposure

Notwithstanding the pop culture images like figure v, it is not possible for a multicellular animal to become “mutated” as a whole. In most cases, a particle of ionizing radiation will not even hit the DNA, and even if it does, it will only affect the DNA of a single cell, not every cell in the animal’s body. Typically, that cell is simply killed, because the DNA becomes unable to function properly. Once in a while, however, the DNA may be altered so as to make that cell cancerous. For instance, skin cancer can be caused by UV light hitting a single skin cell in the body of a sunbather. If that cell becomes cancerous and begins reproducing uncontrollably, she will end up with a tumor twenty years later.

Other than cancer, the only other dramatic effect that can result from altering a single cell’s DNA is if that cell happens to be a sperm or ovum, which can result in nonviable or mutated offspring. Men are relatively immune to reproductive harm from radiation, because their sperm cells are replaced frequently. Women are more vulnerable because they keep the same set of ova as long as they live.

Effects of high doses of radiation

A whole-body exposure of 5,000,000 μSv will kill a person within a week or so. Luckily, only a small number of humans have ever been exposed to such levels: one scientist working on the Manhattan Project, some victims of the Nagasaki and Hiroshima explosions, and 31 workers at Chernobyl. Death occurs by massive killing of cells, especially in the blood-producing cells of the bone marrow.

Effects of low doses radiation

Lower levels, on the order of 1,000,000 μSv , were inflicted on some people at Nagasaki and Hiroshima. No acute symptoms result from this level of exposure, but certain types of cancer are significantly more common among these people. It was originally expected that the radiation would cause many mutations resulting in birth defects, but very few such inherited effects have been observed.

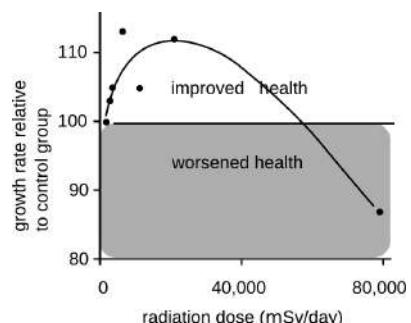
A great deal of time has been spent debating the effects of very low levels of ionizing radiation. The following table gives some sample figures.

maximum <i>beneficial</i> dose per day	$\sim 10,000 \mu\text{Sv}$
CT scan	$\sim 10,000 \mu\text{Sv}$
natural background per year	2,000-7,000 μSv
health guidelines for exposure to a fetus	1,000 μSv
flying from New York to Tokyo	150 μSv
chest x-ray	50 μSv

Note that the largest number, on the first line of the table, is the maximum *beneficial* dose. The most useful evidence comes from

experiments in animals, which can intentionally be exposed to significant and well measured doses of radiation under controlled conditions. Experiments show that low levels of radiation activate cellular damage control mechanisms, increasing the health of the organism. For example, exposure to radiation up to a certain level makes mice grow faster; makes guinea pigs' immune systems function better against diphteria; increases fertility in trout and mice; improves fetal mice's resistance to disease; increases the life-spans of flour beetles and mice; and reduces mortality from cancer in mice. This type of effect is called radiation hormesis.

There is also some evidence that in humans, small doses of radiation increase fertility, reduce genetic abnormalities, and reduce mortality from cancer. The human data, however, tend to be very poor compared to the animal data. Due to ethical issues, one cannot do controlled experiments in humans. For example, one of the best sources of information has been from the survivors of the Hiroshima and Nagasaki bomb blasts, but these people were also exposed to high levels of carcinogenic chemicals in the smoke from their burning cities; for comparison, firefighters have a heightened risk of cancer, and there are also significant concerns about cancer from the 9/11 attacks in New York. The direct empirical evidence about radiation hormesis in humans is therefore not good enough to tell us anything unambiguous,³ and the most scientifically reasonable approach is to assume that the results in animals also hold for humans: small doses of radiation in humans are beneficial, rather than harmful. However, a variety of cultural and historical factors have led to a situation in which public health policy is based on the assumption, known as "linear no-threshold" (LNT), that even tiny doses of radiation are harmful, and that the risk they carry is proportional to the dose. In other words, law and policy are made based on the assumption that the effects of radiation on humans are dramatically different than its effects on mice and guinea pigs. Even with the unrealistic assumption of LNT, one can still evaluate risks by comparing with natural background radiation. For example, we can see that the effect of a chest x-ray is about a hundred times smaller than the effect of spending a year in Colorado, where the level of natural background radiation from cosmic rays is higher than average, due to the high altitude. Dropping the implausible LNT assumption, we can see that the impact on one's health of spending a year in Colorado is likely to be *positive*, because the excess radiation is below the maximum beneficial level.



w / A typical example of radiation hormesis: the health of mice is improved by low levels of radiation. In this study, young mice were exposed to fairly high levels of x-rays, while a control group of mice was not exposed. The mice were weighed, and their rate of growth was taken as a measure of their health. At levels below about 50,000 μSv , the radiation had a beneficial effect on the health of the mice, presumably by activating cellular damage control mechanisms. The two highest data points are statistically significant at the 99% level. The curve is a fit to a theoretical model. Redrawn from T.D. Luckey, *Hormesis with Ionizing Radiation*, CRC Press, 1980.

³For two opposing viewpoints, see Tubiana et al., "The Linear No-Threshold Relationship Is Inconsistent with Radiation Biologic and Experimental Data," Radiology, 251 (2009) 13 and Little et al., "Risks Associated with Low Doses and Low Dose Rates of Ionizing Radiation: Why Linearity May Be (Almost) the Best We Can Do," Radiology, 251 (2009) 6.



x / Wild Przewalski's horses prosper in the Chernobyl area.



y / Fossil fuels have done incomparably more damage to the environment than nuclear power ever has. Polar bears' habitat is rapidly being destroyed by global warming.

The green case for nuclear power

In the late twentieth century, antinuclear activists largely succeeded in bringing construction of new nuclear power plants to a halt in the U.S. Ironically, we now know that the burning of fossil fuels, which leads to global warming, is a far more grave threat to the environment than even the Chernobyl disaster. A team of biologists writes: “During recent visits to Chernobyl, we experienced numerous sightings of moose (*Alces alces*), roe deer (*Capreolus capreolus*), Russian wild boar (*Sus scrofa*), foxes (*Vulpes vulpes*), river otter (*Lutra canadensis*), and rabbits (*Lepus europaeus*) ... Diversity of flowers and other plants in the highly radioactive regions is impressive and equals that observed in protected habitats outside the zone ... The observation that typical human activity (industrialization, farming, cattle raising, collection of firewood, hunting, etc.) is more devastating to biodiversity and abundance of local flora and fauna than is the worst nuclear power plant disaster validates the negative impact the exponential growth of human populations has on wildlife.”⁴

Nuclear power is the only source of energy that is sufficient to replace any significant percentage of energy from fossil fuels on the rapid schedule demanded by the speed at which global warming is progressing. People worried about the downside of nuclear energy might be better off putting their energy into issues related to nuclear weapons: the poor stewardship of the former Soviet Union’s warheads; nuclear proliferation in unstable states such as Pakistan; and the poor safety and environmental history of the superpowers’ nuclear weapons programs, including the loss of several warheads in plane crashes, and the environmental disaster at the Hanford, Washington, weapons plant.

Protection from radiation

People do sometimes work with strong enough radioactivity that there is a serious health risk. Typically the scariest sources are those used in cancer treatment and in medical and biological research. Also, a dental technician, for example, needs to take precautions to avoid accumulating a large radiation dose from giving dental x-rays to many patients. There are three general ways to reduce exposure: time, distance, and shielding. This is why a dental technician doing x-rays wears a lead apron (shielding) and steps outside of the x-ray room while running an exposure (distance). Reducing the time of exposure dictates, for example, that a person working with a hot cancer-therapy source would minimize the amount of time spent

⁴Baker and Chesser, Env. Toxicology and Chem. 19 (1231) 2000. Similar effects have been seen at the Bikini Atoll, the site of a 1954 hydrogen bomb test. Although some species have disappeared from the area, the coral reef is in many ways healthier than similar reefs elsewhere, because humans have tended to stay away for fear of radiation (Richards et al., Marine Pollution Bulletin 56 (2008) 503).

near it.

Shielding against alpha and beta particles is trivial to accomplish. (Alphas can't even penetrate the skin.) Gammas and x-rays interact most strongly with materials that are dense and have high atomic numbers, which is why lead is so commonly used. But other materials will also work. For example, the reason that bones show up so clearly on x-ray images is that they are dense and contain plenty of calcium, which has a higher atomic number than the elements found in most other body tissues, which are mostly made of water.

Neutrons are difficult to shield against. Because they are electrically neutral, they don't interact intensely with matter in the same way as alphas and betas. They only interact if they happen to collide head-on with a nucleus, and that doesn't happen very often because nuclei are tiny targets. Kinematically, a collision can transfer kinetic energy most efficiently when the target is as low in mass as possible compared to the projectile. For this reason, substances that contain a lot of hydrogen make the best shielding against neutrons. Blocks of paraffin wax from the supermarket are often used for this purpose.

8.2.10 * The creation of the elements

Creation of hydrogen and helium in the Big Bang

Did all the chemical elements we're made of come into being in the big bang?⁵ Temperatures in the first microseconds after the big bang were so high that atoms and nuclei could not hold together at all. After things had cooled down enough for nuclei and atoms to exist, there was a period of about three minutes during which the temperature and density were high enough for fusion to occur, but not so high that atoms could hold together. We have a good, detailed understanding of the laws of physics that apply under these conditions, so theorists are able to say with confidence that the only element heavier than hydrogen that was created in significant quantities was helium.

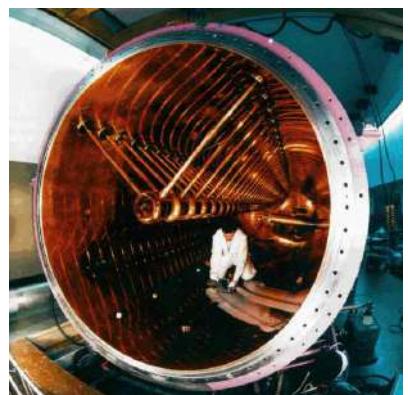
We are stardust

In that case, where did all the other elements come from? Astronomers came up with the answer. By studying the combinations of wavelengths of light, called spectra, emitted by various stars, they had been able to determine what kinds of atoms they contained. (We will have more to say about spectra at the end of this book.) They found that the stars fell into two groups. One type was nearly 100% hydrogen and helium, while the other contained 99% hydrogen and helium and 1% other elements. They interpreted these as two

⁵The evidence for the big bang theory of the origin of the universe was discussed on p. 370.



z / The Crab Nebula is a remnant of a supernova explosion. Almost all the elements our planet is made of originated in such explosions.



aa / Construction of the UNILAC accelerator in Germany, one of whose uses is for experiments to create very heavy artificial elements. In such an experiment, fusion products recoil through a device called SHIP (not shown) that separates them based on their charge-to-mass ratios — it is essentially just a scaled-up version of Thomson's apparatus. A typical experiment runs for several months, and out of the billions of fusion reactions induced during this time, only one or two may result in the production of superheavy atoms. In all the rest, the fused nucleus breaks up immediately. SHIP is used to identify the small number of “good” reactions and separate them from this intense background.

generations of stars. The first generation had formed out of clouds of gas that came fresh from the big bang, and their composition reflected that of the early universe. The nuclear fusion reactions by which they shine have mainly just increased the proportion of helium relative to hydrogen, without making any heavier elements. The members of the first generation that we see today, however, are only those that lived a long time. Small stars are more miserly with their fuel than large stars, which have short lives. The large stars of the first generation have already finished their lives. Near the end of its lifetime, a star runs out of hydrogen fuel and undergoes a series of violent and spectacular reorganizations as it fuses heavier and heavier elements. Very large stars finish this sequence of events by undergoing supernova explosions, in which some of their material is flung off into the void while the rest collapses into an exotic object such as a black hole or neutron star.

The second generation of stars, of which our own sun is an example, condensed out of clouds of gas that had been enriched in heavy elements due to supernova explosions. It is those heavy elements that make up our planet and our bodies.

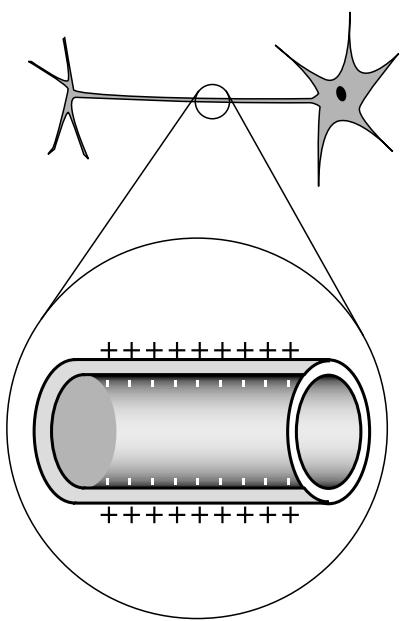
Discussion Questions

- A** Should the quality factor for neutrinos be very small, because they mostly don't interact with your body?
- B** Would an alpha source be likely to cause different types of cancer depending on whether the source was external to the body or swallowed in contaminated food? What about a gamma source?

Problems

The symbols \checkmark , \blacksquare , etc. are explained on page 528.

1 The figure shows a neuron, which is the type of cell your nerves are made of. Neurons serve to transmit sensory information to the brain, and commands from the brain to the muscles. All this data is transmitted electrically, but even when the cell is resting and not transmitting any information, there is a layer of negative electrical charge on the inside of the cell membrane, and a layer of positive charge just outside it. This charge is in the form of various ions dissolved in the interior and exterior fluids. Why would the negative charge remain plastered against the inside surface of the membrane, and likewise why doesn't the positive charge wander away from the outside surface? \blacksquare



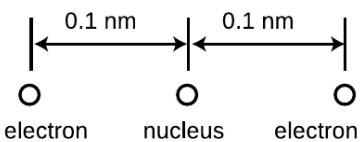
2 The Earth and Moon are bound together by gravity. If, instead, the force of attraction were the result of each having a charge of the same magnitude but opposite in sign, find the quantity of charge that would have to be placed on each to produce the required force. $\checkmark \blacksquare$

3 A helium atom finds itself momentarily in this arrangement. Find the direction and magnitude of the force acting on the right-hand electron. The two protons in the nucleus are so close together (~ 1 fm) that you can consider them as being right on top of each other. $\checkmark \blacksquare$

4 ^{241}Pu decays either by electron decay or by alpha decay. (A given ^{241}Pu nucleus may do either one; it's random.) What are the isotopes created as products of these two modes of decay? \blacksquare

5 Suppose that a proton in a lead nucleus wanders out to the surface of the nucleus, and experiences a strong nuclear force of about 8 kN from the nearby neutrons and protons pulling it back in. Compare this numerically to the repulsive electrical force from the other protons, and verify that the net force is attractive. A lead nucleus is very nearly spherical, is about 6.5 fm in radius, and contains 82 protons, each with a charge of $+e$, where $e = 1.60 \times 10^{-19} \text{ C}$. $\checkmark \blacksquare$

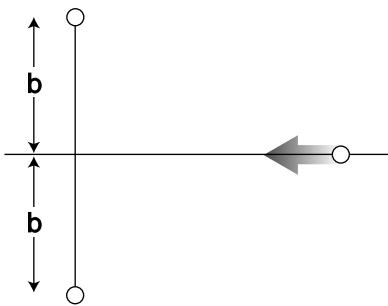
Problem 1. Top: A realistic picture of a neuron. Bottom: A simplified diagram of one segment of the tail (axon).



Problem 3.

6 The nuclear process of beta decay by electron capture is described parenthetically on page 512. The reaction is $\text{p} + \text{e}^- \rightarrow \text{n} + \nu$.

- Show that charge is conserved in this reaction.
- Conversion between energy and mass is discussed in chapter 7. Based on these ideas, explain why electron capture doesn't occur in hydrogen atoms. (If it did, matter wouldn't exist!) \blacksquare

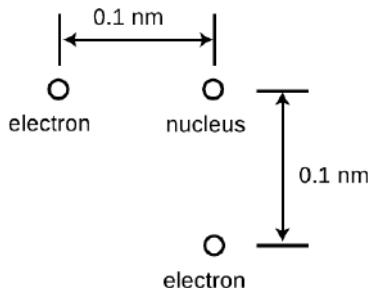


Problem 7.

7 In the semifinals of an electrostatic croquet tournament, Jessica hits her positively charged ball, sending it across the playing field, rolling to the left along the x axis. It is repelled by two other positive charges. These two equal charges are fixed on the y axis at the locations shown in the figure. (a) Express the force on the ball in terms of the ball's position, x . (b) At what value of x does the ball experience the greatest deceleration? Express your answer in terms of b . [Based on a problem by Halliday and Resnick.] ✓ ■

8 Suppose that at some instant in time, a wire extending from $x = 0$ to $x = \infty$ holds a charge density, in units of coulombs per meter, given by ae^{-bx} . This type of charge density, dq/dx , is typically denoted as λ (Greek letter lambda). Find the total charge on the wire. ✓ ■

9 Use the nutritional information on some packaged food to make an order-of-magnitude estimate of the amount of chemical energy stored in one atom of food, in units of joules. Assume that a typical atom has a mass of 10^{-26} kg. This constitutes a rough estimate of the amounts of energy there are on the atomic scale. [See chapter 0 for help on how to do order-of-magnitude estimates. Note that a nutritional “calorie” is really a kilocalorie; see page 1072.] ✓ ■



Problem 10.

10 The helium atom of problem 3 has some new experiences, goes through some life changes, and later on finds itself in the configuration shown here. What are the direction and magnitude of the force acting on the bottom electron? (Draw a sketch to make clear the definition you are using for the angle that gives direction.) ✓ ■

11 A neon light consists of a long glass tube full of neon, with metal caps on the ends. Positive charge is placed on one end of the tube, negative on the other. The electric forces generated can be strong enough to strip electrons off of a certain number of neon atoms. Assume for simplicity that only one electron is ever stripped off of any neon atom. When an electron is stripped off of an atom, both the electron and the neon atom (now an ion) have electric charge, and they are accelerated by the forces exerted by the charged ends of the tube. (They do not feel any significant forces from the other ions and electrons within the tube, because only a tiny minority of neon atoms ever gets ionized.) Light is finally produced when ions are reunited with electrons. Give a numerical comparison of the magnitudes and directions of the accelerations of the electrons and ions. [You may need some data from appendix 4, p. 1070.] ✓ ■

12 The subatomic particles called muons behave exactly like electrons, except that a muon's mass is greater by a factor of 206.77. Muons are continually bombarding the Earth as part of the stream of particles from space known as cosmic rays. When a muon strikes an atom, it can displace one of its electrons. If the atom happens to be a hydrogen atom, then the muon takes up an orbit that is on the average 206.77 times closer to the proton than the orbit of the ejected electron. How many times greater is the electric force experienced by the muon than that previously felt by the electron?

✓ ■

13 (a) Recall that the gravitational energy of two gravitationally interacting spheres is given by $U_g = -Gm_1m_2/r$, where r is the center-to-center distance. What would be the analogous equation for two electrically interacting spheres? Justify your choice of a plus or minus sign on physical grounds, considering attraction and repulsion.

✓

(b) Use this expression to estimate the energy required to pull apart a raisin-cookie atom of the one-electron type, assuming a radius of 10^{-10} m.

✓

(c) Compare this with the result of problem 9.

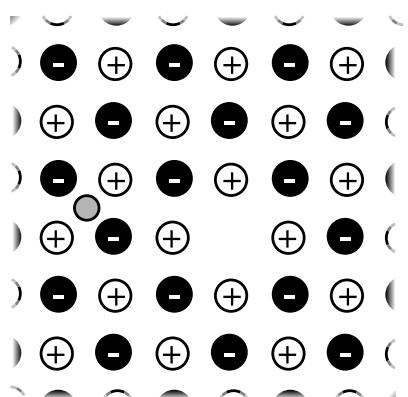
■

14 If you put two hydrogen atoms near each other, they will feel an attractive force, and they will pull together to form a molecule. (Molecules consisting of two hydrogen atoms are the normal form of hydrogen gas.) Why do they feel a force if they are near each other, since each is electrically neutral? Shouldn't the attractive and repulsive forces all cancel out exactly? Use the raisin cookie model. (Students who have taken chemistry often try to use fancier models to explain this, but if you can't explain it using a simple model, you probably don't understand the fancy model as well as you thought you did!)

■

15 The figure shows one layer of the three-dimensional structure of a salt crystal. The atoms extend much farther off in all directions, but only a six-by-six square is shown here. The larger circles are the chlorine ions, which have charges of $-e$, where $e = 1.60 \times 10^{-19}$ C. The smaller circles are sodium ions, with charges of $+e$. The center-to-center distance between neighboring ions is about 0.3 nm. Real crystals are never perfect, and the crystal shown here has two defects: a missing atom at one location, and an extra lithium atom, shown as a grey circle, inserted in one of the small gaps. If the lithium atom has a charge of $+e$, what is the direction and magnitude of the total force on it? Assume there are no other defects nearby in the crystal besides the two shown here.

▷ Hint, p. 1036 ✓ ■



Problem 15.

16 Suppose you are holding your hands in front of you, 10 cm apart.

- (a) Estimate the total number of electrons in each hand. ✓
- (b) Estimate the total repulsive force of all the electrons in one hand on all the electrons in the other. ✓
- (c) Why don't you feel your hands repelling each other?
- (d) Estimate how much the charge of a proton could differ in magnitude from the charge of an electron without creating a noticeable force between your hands. ■

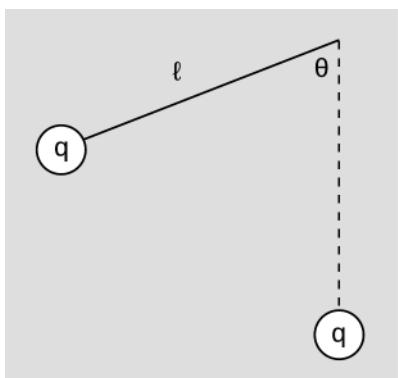
17 Potassium 40 is the strongest source of naturally occurring beta radioactivity in our environment. It decays according to



The energy released in the decay is 1.33 MeV, where 1 eV is defined as the fundamental charge e multiplied by one volt. The energy is shared randomly among the products, subject to the constraint imposed by conservation of energy-momentum, which dictates that very little of the energy is carried by the recoiling calcium nucleus. Determine the maximum energy of the calcium, and compare with the typical energy of a chemical bond, which is a few eV. If the potassium is part of a molecule, do we expect the molecule to survive? Carry out the calculation first by assuming that the electron is ultrarelativistic, then without the approximation, and comment on the how good the approximation is. ■

18 Several pointlike, interacting particles are released, all initially at rest, within the same finite region of space. We want to know whether, without violating conservation of energy, it is possible for one of these particles to be ejected to an infinite distance, while the others stay within the original region. Consider this question for each of the following systems. (a) Two particles of mass m , interacting gravitationally. (b) Three such particles. (c) Two particles, both with charge q . (d) Charges q and $-q$. (e) Charges q , q , and $-q$. ■

19 As shown in the figure, a particle of mass m and charge q hangs from a string of length ℓ , forming a pendulum fixed at a central point. Another charge q is fixed at the same distance ℓ , directly below the center. Find the equilibrium values of θ and determine whether they are stable or unstable. ■



Problem 19.

Key to symbols:

■ easy ■ typical ■ challenging ■ difficult ■ very difficult

✓ An answer check is available at www.lightandmatter.com.

Exercises

Exercise 8A: Nuclear decay

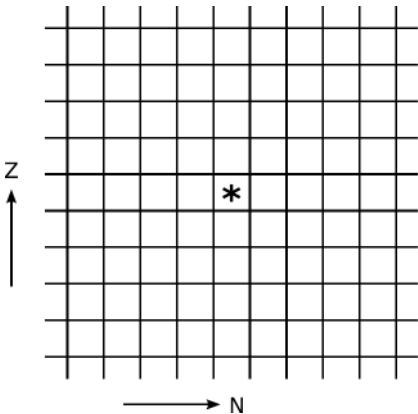
1. Consulting a periodic table, find the N , Z , and A of the following:

	N	Z	A
${}^4\text{He}$			
${}^{244}\text{Pu}$			

2. Consider the following five decay processes:

- α decay
- γ decay
- $\text{p} \rightarrow \text{n} + \text{e}^+ + \nu$ (β^+ decay)
- $\text{n} \rightarrow \text{p} + \text{e}^- + \bar{\nu}$ (β^- decay)
- $\text{p} + \text{e}^- \rightarrow \text{n} + \nu$ (electron capture)

What would be the action of each of these on the chart of the nuclei? The * represents the original nucleus.



3. (a) Suppose that ${}^{244}\text{Pu}$ undergoes perfectly symmetric fission, and also emits two neutrons. Find the daughter isotope.

(b) Is the daughter stable, or is it neutron-rich or -poor relative to the line of stability? (To estimate what's stable, you can use a large chart of the nuclei, or, if you don't have one handy, consult a periodic table and use the average atomic mass as an approximation to the stable value of A .)

(c) Consulting the chart of the nuclei (fig. s on p. 518), explain why it turns out this way.

(d) If the daughter is unstable, which process from question #2 would you expect it to decay by?



Chapter 9

Circuits

Madam, what good is a baby? *Michael Faraday, when asked by Queen Victoria what the electrical devices in his lab were good for*

A few years ago, my wife and I bought a house with Character, Character being a survival mechanism that houses have evolved in order to convince humans to agree to much larger mortgage payments than they'd originally envisioned. Anyway, one of the features that gives our house Character is that it possesses, built into the wall of the family room, a set of three pachinko machines. These are Japanese gambling devices sort of like vertical pinball machines. (The legal papers we got from the sellers hastened to tell us that they were "for amusement purposes only.") Unfortunately, only one of the three machines was working when we moved in, and it soon died on us. Having become a pachinko addict, I decided to fix it, but that was easier said than done. The inside is a veritable Rube Goldberg mechanism of levers, hooks, springs, and chutes. My hormonal pride, combined with my Ph.D. in physics, made me certain of success, and rendered my eventual utter failure all the more demoralizing.

Contemplating my defeat, I realized how few complex mechanical devices I used from day to day. Apart from our cars and my

saxophone, every technological tool in our modern life-support system was electronic rather than mechanical.

9.1 Current and voltage

9.1.1 Current

Unity of all types of electricity

We are surrounded by things we have been *told* are “electrical,” but it’s far from obvious what they have in common to justify being grouped together. What relationship is there between the way socks cling together and the way a battery lights a lightbulb? We have been told that both an electric eel and our own brains are somehow electrical in nature, but what do they have in common?

British physicist Michael Faraday (1791-1867) set out to address this problem. He investigated electricity from a variety of sources — including electric eels! — to see whether they could all produce the same effects, such as shocks and sparks, attraction and repulsion. “Heating” refers, for example, to the way a lightbulb filament gets hot enough to glow and emit light. Magnetic induction is an effect discovered by Faraday himself that connects electricity and magnetism. We will not study this effect, which is the basis for the electric generator, in detail until later in the book.

	attraction and			
	shocks	sparks	repulsion	heating
rubbing	✓	✓	✓	✓
battery	✓	✓	✓	✓
animal	✓	✓	(✓)	✓
magnetically induced	✓	✓	✓	✓

The table shows a summary of some of Faraday’s results. Check marks indicate that Faraday or his close contemporaries were able to verify that a particular source of electricity was capable of producing a certain effect. (They evidently failed to demonstrate attraction and repulsion between objects charged by electric eels, although modern workers have studied these species in detail and been able to understand all their electrical characteristics on the same footing as other forms of electricity.)

Faraday’s results indicate that there is nothing fundamentally different about the types of electricity supplied by the various sources. They are all able to produce a wide variety of identical effects. Wrote Faraday, “The general conclusion which must be drawn from this collection of facts is that electricity, whatever may be its source, is identical in its nature.”

If the types of electricity are the same thing, what thing is that? The answer is provided by the fact that all the sources of electricity can cause objects to repel or attract each other. We use the word



a / *Gymnotus carapo*, a knifefish, uses electrical signals to sense its environment and to communicate with others of its species.
(Greg DeGreef)

“charge” to describe the property of an object that allows it to participate in such electrical forces, and we have learned that charge is present in matter in the form of nuclei and electrons. Evidently all these electrical phenomena boil down to the motion of charged particles in matter.

Electric current

If the fundamental phenomenon is the motion of charged particles, then how can we define a useful numerical measurement of it? We might describe the flow of a river simply by the velocity of the water, but velocity will not be appropriate for electrical purposes because we need to take into account how much charge the moving particles have, and in any case there are no practical devices sold at Radio Shack that can tell us the velocity of charged particles. Experiments show that the intensity of various electrical effects is related to a different quantity: the number of coulombs of charge that pass by a certain point per second. By analogy with the flow of water, this quantity is called the electric *current*, I . Its units of coulombs/second are more conveniently abbreviated as amperes, 1 A=1 C/s. (In informal speech, one usually says “amps.”)

The main subtlety involved in this definition is how to account for the two types of charge. The stream of water coming from a hose is made of atoms containing charged particles, but it produces none of the effects we associate with electric currents. For example, you do not get an electrical shock when you are sprayed by a hose. This type of experiment shows that the effect created by the motion of one type of charged particle can be canceled out by the motion of the opposite type of charge in the same direction. In water, every oxygen atom with a charge of $+8e$ is surrounded by eight electrons with charges of $-e$, and likewise for the hydrogen atoms.

We therefore refine our definition of current as follows:

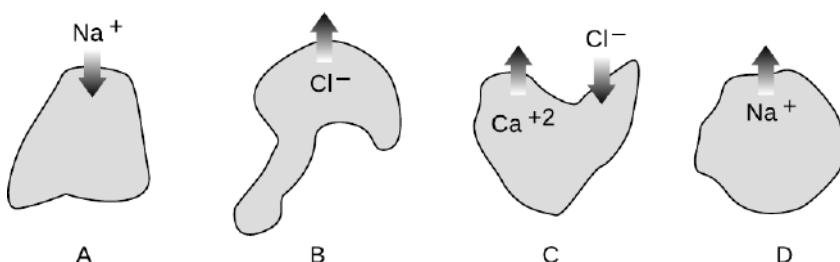
When charged particles are exchanged between regions of space A and B, the electric current flowing from A to B is defined as

$$I = \frac{dq}{dt},$$

where dq is the change in region B’s total charge occurring over a period of time dt .

In the garden hose example, your body picks up equal amounts of positive and negative charge, resulting in no change in your total charge, so the electrical current flowing into you is zero.

b / Example 1



Ions moving across a cell membrane

example 1

▷ Figure b shows ions, labeled with their charges, moving in or out through the membranes of four cells. If the ions all cross the membranes during the same interval of time, how would the currents into the cells compare with each other?

▷ We're just assuming the rate of flow is constant, so we can talk about Δq instead of dq .

Cell A has positive current going into it because its charge is increased, i.e., has a positive value of Δq .

Cell B has the same current as cell A, because by losing one unit of negative charge it also ends up increasing its own total charge by one unit.

Cell C's total charge is reduced by three units, so it has a large negative current going into it.

Cell D loses one unit of charge, so it has a small negative current into it.

Finding current given charge

example 2

▷ A charged balloon falls to the ground, and its charge begins leaking off to the Earth. Suppose that the charge on the balloon is given by $q = ae^{-bt}$. Find the current as a function of time, and interpret the answer.

▷ Taking the derivative, we have

$$I = \frac{dq}{dt} \\ = -abe^{-bt}$$

An exponential function approaches zero as the exponent gets more and more negative. This means that both the charge and the current are decreasing in magnitude with time. It makes sense that the charge approaches zero, since the balloon is losing its charge. It also makes sense that the current is decreasing in magnitude, since charge cannot flow at the same rate forever without overshooting zero.

The reverse of differentiation is integration, so if we know the

current as a function of time, we can find the charge by integrating. Example 8 on page 544 shows such a calculation.

It may seem strange to say that a negatively charged particle going one way creates a current going the other way, but this is quite ordinary. As we will see, currents flow through metal wires via the motion of electrons, which are negatively charged, so the direction of motion of the electrons in a circuit is always opposite to the direction of the current. Of course it would have been convenient if Benjamin Franklin had defined the positive and negative signs of charge the opposite way, since so many electrical devices are based on metal wires.

Number of electrons flowing through a lightbulb example 3

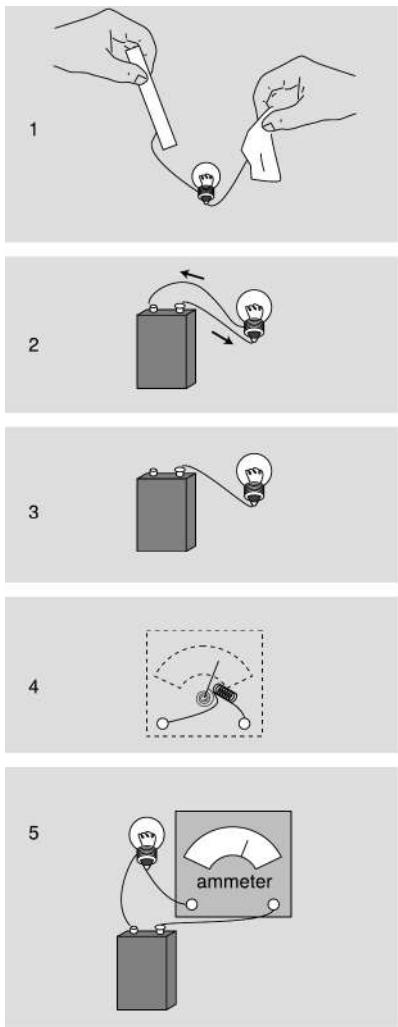
- ▷ If a lightbulb has 1.0 A flowing through it, how many electrons will pass through the filament in 1.0 s?
- ▷ We are only calculating the number of electrons that flow, so we can ignore the positive and negative signs. Also, since the rate of flow is constant, we don't really need to think in terms of calculus; the derivative dq/dt that defines current is the same as $\Delta q/\Delta t$ in this situation. Solving for $\Delta q = I\Delta t$ gives a charge of 1.0 C flowing in this time interval. The number of electrons is

$$\begin{aligned}\text{number of electrons} &= \text{coulombs} \times \frac{\text{electrons}}{\text{coulomb}} \\ &= \text{coulombs} / \frac{\text{coulombs}}{\text{electron}} \\ &= 1.0 \text{ C}/e \\ &= 6.2 \times 10^{18}\end{aligned}$$

9.1.2 Circuits

How can we put electric currents to work? The only method of controlling electric charge we have studied so far is to charge different substances, e.g., rubber and fur, by rubbing them against each other. Figure c/1 shows an attempt to use this technique to light a lightbulb. This method is unsatisfactory. True, current will flow through the bulb, since electrons can move through metal wires, and the excess electrons on the rubber rod will therefore come through the wires and bulb due to the attraction of the positively charged fur and the repulsion of the other electrons. The problem is that after a zillionth of a second of current, the rod and fur will both have run out of charge. No more current will flow, and the lightbulb will go out.

Figure c/2 shows a setup that works. The battery pushes charge through the circuit, and recycles it over and over again. (We will have more to say later in this chapter about how batteries work.) This is called a *complete circuit*. Today, the electrical use of the word “circuit” is the only one that springs to mind for most people,



c / 1. Static electricity runs out quickly. 2. A practical circuit. 3. An open circuit. 4. How an ammeter works. 5. Measuring the current with an ammeter.

but the original meaning was to travel around and make a round trip, as when a circuit court judge would ride around the boondocks, dispensing justice in each town on a certain date.

Note that an example like c/3 does not work. The wire will quickly begin acquiring a net charge, because it has no way to get rid of the charge flowing into it. The repulsion of this charge will make it more and more difficult to send any more charge in, and soon the electrical forces exerted by the battery will be canceled out completely. The whole process would be over so quickly that the filament would not even have enough time to get hot and glow. This is known as an *open circuit*. Exactly the same thing would happen if the complete circuit of figure c/2 was cut somewhere with a pair of scissors, and in fact that is essentially how an ordinary light switch works: by opening up a gap in the circuit.

The definition of electric current we have developed has the great virtue that it is easy to measure. In practical electrical work, one almost always measures current, not charge. The instrument used to measure current is called an *ammeter*. A simplified ammeter, c/4, simply consists of a coiled-wire magnet whose force twists an iron needle against the resistance of a spring. The greater the current, the greater the force. Although the construction of ammeters may differ, their use is always the same. We break into the path of the electric current and interpose the meter like a tollbooth on a road, c/5. There is still a complete circuit, and as far as the battery and bulb are concerned, the ammeter is just another segment of wire.

Does it matter where in the circuit we place the ammeter? Could we, for instance, have put it in the left side of the circuit instead of the right? Conservation of charge tells us that this can make no difference. Charge is not destroyed or “used up” by the lightbulb, so we will get the same current reading on either side of it. What is “used up” is energy stored in the battery, which is being converted into heat and light energy.

9.1.3 Voltage

The volt unit

Electrical circuits can be used for sending signals, storing information, or doing calculations, but their most common purpose by far is to manipulate energy, as in the battery-and-bulb example of the previous section. We know that lightbulbs are rated in units of watts, i.e., how many joules per second of energy they can convert into heat and light, but how would this relate to the flow of charge as measured in amperes? By way of analogy, suppose your friend, who didn’t take physics, can’t find any job better than pitching bales of hay. The number of calories he burns per hour will certainly depend on how many bales he pitches per minute, but it will also be proportional to how much mechanical work he has to do on each bale.

If his job is to toss them up into a hayloft, he will get tired a lot more quickly than someone who merely tips bales off a loading dock into trucks. In metric units,

$$\frac{\text{joules}}{\text{second}} = \frac{\text{haybales}}{\text{second}} \times \frac{\text{joules}}{\text{haybale}}.$$

Similarly, the rate of energy transformation by a battery will not just depend on how many coulombs per second it pushes through a circuit but also on how much mechanical work it has to do on each coulomb of charge:

$$\frac{\text{joules}}{\text{second}} = \frac{\text{coulombs}}{\text{second}} \times \frac{\text{joules}}{\text{coulomb}}$$

or

$$\text{power} = \text{current} \times \text{work per unit charge}.$$

Units of joules per coulomb are abbreviated as *volts*, $1 \text{ V}=1 \text{ J/C}$, named after the Italian physicist Alessandro Volta. Everyone knows that batteries are rated in units of volts, but the voltage concept is more general than that; it turns out that voltage is a property of every point in space. To gain more insight, let's think more carefully about what goes on in the battery and bulb circuit.

The concept of voltage (electrical potential) in general

To do work on a charged particle, the battery apparently must be exerting forces on it. How does it do this? Well, the only thing that can exert an electrical force on a charged particle is another charged particle. It's as though the haybales were pushing and pulling each other into the hayloft! This is potentially a horribly complicated situation. Even if we knew how much excess positive or negative charge there was at every point in the circuit (which realistically we don't) we would have to calculate zillions of forces using Coulomb's law, perform all the vector additions, and finally calculate how much work was being done on the charges as they moved along. To make things even more scary, there is more than one type of charged particle that moves: electrons are what move in the wires and the bulb's filament, but ions are the moving charge carriers inside the battery. Luckily, there are two ways in which we can simplify things:

The situation is unchanging. Unlike the imaginary setup in which we attempted to light a bulb using a rubber rod and a piece of fur, this circuit maintains itself in a steady state (after perhaps a microsecond-long period of settling down after the circuit is first assembled). The current is steady, and as charge flows out of any area of the circuit it is replaced by the same amount of charge flowing in. The amount of excess positive or negative charge in any part of the circuit therefore stays constant. Similarly, when we watch a river flowing, the water goes by but the river doesn't disappear.

Force depends only on position. Since the charge distribution is not changing, the total electrical force on a charged particle depends only on its own charge and on its location. If another charged particle of the same type visits the same location later on, it will feel exactly the same force.

The second observation tells us that there is nothing all that different about the experience of one charged particle as compared to another's. If we single out one particle to pay attention to, and figure out the amount of work done on it by electrical forces as it goes from point A to point B along a certain path, then this is the same amount of work that will be done on any other charged particles of the same type as it follows the same path. For the sake of visualization, let's think about the path that starts at one terminal of the battery, goes through the light bulb's filament, and ends at the other terminal. When an object experiences a force that depends only on its position (and when certain other, technical conditions are satisfied), we can define an electrical energy associated with the position of that object. The amount of work done on the particle by electrical forces as it moves from A to B equals the drop in electrical energy between A and B. This electrical energy is what is being converted into other forms of energy such as heat and light. We therefore define ΔV in general as electrical energy per unit charge:

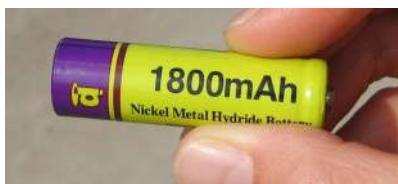
The ΔV between two points in space is defined as

$$\Delta V = \Delta U_{elec}/q,$$

where ΔU_{elec} is the change in the electrical energy of a particle with charge q as it moves from the initial point to the final point. In this context, where we think of the voltage as being a scalar function that is defined everywhere in space, it is more common in formal writing to refer to it as the electrical *potential*.

The amount of power dissipated (i.e., rate at which energy is transformed by the flow of electricity) is then given by the equation

$$P = I\Delta V.$$



d / Example 4.

Energy stored in a battery

example 4

- ▷ The 1.2 V rechargeable battery in figure d is labeled 1800 milliamp-hours. What is the maximum amount of energy the battery can store?
- ▷ An ampere-hour is a unit of current multiplied by a unit of time. Current is charge per unit time, so an ampere-hour is in fact a funny unit of *charge*:

$$(1 \text{ A})(1 \text{ hour}) = (1 \text{ C/s})(3600 \text{ s}) \\ = 3600 \text{ C}$$

1800 milliamp-hours is therefore $1800 \times 10^{-3} \times 3600 \text{ C} = 6.5 \times 10^3 \text{ C}$. That's a huge number of charged particles, but the total loss of electrical energy will just be their total charge multiplied by the voltage difference across which they move:

$$\begin{aligned}\Delta U_{\text{elec}} &= q\Delta V \\ &= (6.5 \times 10^3 \text{ C})(1.2 \text{ V}) \\ &= 7.8 \text{ kJ}\end{aligned}$$

Units of volt-amps

example 5

- ▷ Doorbells are often rated in volt-amps. What does this combination of units mean?
- ▷ Current times voltage gives units of power, $P = I\Delta V$, so volt-amps are really just a nonstandard way of writing watts. They are telling you how much power the doorbell requires.

Power dissipated by a battery and bulb

example 6

- ▷ If a 9.0-volt battery causes 1.0 A to flow through a lightbulb, how much power is dissipated?
- ▷ The voltage rating of a battery tells us what voltage difference ΔV it is designed to maintain between its terminals.

$$\begin{aligned}P &= I \Delta V \\ &= 9.0 \text{ A} \cdot \text{V} \\ &= 9.0 \frac{\text{C}}{\text{s}} \cdot \frac{\text{J}}{\text{C}} \\ &= 9.0 \text{ J/s} \\ &= 9.0 \text{ W}\end{aligned}$$

The only nontrivial thing in this problem was dealing with the units. One quickly gets used to translating common combinations like $\text{A} \cdot \text{V}$ into simpler terms.

Here are a few questions and answers about the voltage concept.

Question: OK, so what *is* voltage, really?

Answer: A device like a battery has positive and negative charges inside it that push other charges around the outside circuit. A higher-voltage battery has denser charges in it, which will do more work on each charged particle that moves through the outside circuit.

To use a gravitational analogy, we can put a paddlewheel at the bottom of either a tall waterfall or a short one, but a kg of water that falls through the greater gravitational energy difference will have more energy to give up to the paddlewheel at the bottom.

Question: Why do we define voltage as electrical energy divided by charge, instead of just defining it as electrical energy?

Answer: One answer is that it's the only definition that makes the equation $P = I\Delta V$ work. A more general answer is that we want to be able to define a voltage difference between any two points in space without having to know in advance how much charge the particles moving between them will have. If you put a nine-volt battery on your tongue, then the charged particles that move across your tongue and give you that tingly sensation are not electrons but ions, which may have charges of $+e$, $-2e$, or practically anything. The manufacturer probably expected the battery to be used mostly in circuits with metal wires, where the charged particles that flowed would be electrons with charges of $-e$. If the ones flowing across your tongue happen to have charges of $-2e$, the electrical energy difference for them will be twice as much, but dividing by their charge of $-2e$ in the definition of voltage will still give a result of 9 V.

Question: Are there two separate roles for the charged particles in the circuit, a type that sits still and exerts the forces, and another that moves under the influence of those forces?

Answer: No. Every charged particle simultaneously plays both roles. Newton's third law says that any particle that has an electrical force acting on it must also be exerting an electrical force back on the other particle. There are no "designated movers" or "designated force-makers."

Question: Why does the definition of voltage only refer to voltage differences?

Answer: It's perfectly OK to define voltage as $V = U_{elec}/q$. But recall that it is only *differences* in interaction energy, U , that have direct physical meaning in physics. Similarly, voltage differences are really more useful than absolute voltages. A voltmeter measures voltage differences, not absolute voltages.

Discussion Questions

A A roller coaster is sort of like an electric circuit, but it uses gravitational forces on the cars instead of electric ones. What would a high-voltage roller coaster be like? What would a high-current roller coaster be like?

B Criticize the following statements:

“He touched the wire, and 10000 volts went through him.”

“That battery has a charge of 9 volts.”

“You used up the charge of the battery.”

C When you touch a 9-volt battery to your tongue, both positive and negative ions move through your saliva. Which ions go which way?

D I once touched a piece of physics apparatus that had been wired incorrectly, and got a several-thousand-volt voltage difference across my hand. I was not injured. For what possible reason would the shock have had insufficient power to hurt me?

9.1.4 Resistance

Resistance

So far we have simply presented it as an observed fact that a battery-and-bulb circuit quickly settles down to a steady flow, but why should it? Newton’s second law, $a = F/m$, would seem to predict that the steady forces on the charged particles should make them whip around the circuit faster and faster. The answer is that as charged particles move through matter, there are always forces, analogous to frictional forces, that resist the motion. These forces need to be included in Newton’s second law, which is really $a = F_{total}/m$, not $a = F/m$. If, by analogy, you push a crate across the floor at constant speed, i.e., with zero acceleration, the total force on it must be zero. After you get the crate going, the floor’s frictional force is exactly canceling out your force. The chemical energy stored in your body is being transformed into heat in the crate and the floor, and no longer into an increase in the crate’s kinetic energy. Similarly, the battery’s internal chemical energy is converted into heat, not into perpetually increasing the charged particles’ kinetic energy. Changing energy into heat may be a nuisance in some circuits, such as a computer chip, but it is vital in an incandescent lightbulb, which must get hot enough to glow. Whether we like it or not, this kind of heating effect is going to occur any time charged particles move through matter.

What determines the amount of heating? One flashlight bulb designed to work with a 9-volt battery might be labeled 1.0 watts, another 5.0. How does this work? Even without knowing the details of this type of friction at the atomic level, we can relate the heat dissipation to the amount of current that flows via the equation $P = I\Delta V$. If the two flashlight bulbs can have two different values of P when used with a battery that maintains the same ΔV , it



e / Georg Simon Ohm (1787-1854).

must be that the 5.0-watt bulb allows five times more current to flow through it.

For many substances, including the tungsten from which lightbulb filaments are made, experiments show that the amount of current that will flow through it is directly proportional to the voltage difference placed across it. For an object made of such a substance, we define its electrical *resistance* as follows:

If an object inserted in a circuit displays a current flow which is proportional to the voltage difference across it, then we define its resistance as the constant ratio

$$R = \Delta V/I.$$

The units of resistance are volts/ampere, usually abbreviated as ohms, symbolized with the capital Greek letter omega, Ω .

Resistance of a lightbulb

example 7

▷ A flashlight bulb powered by a 9-volt battery has a resistance of $10\ \Omega$. How much current will it draw?

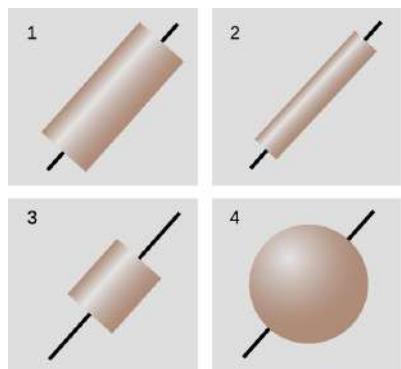
▷ Solving the definition of resistance for I , we find

$$\begin{aligned} I &= \Delta V/R \\ &= 0.9\ \text{V}/\Omega \\ &= 0.9\ \text{V}/(\text{V/A}) \\ &= 0.9\ \text{A} \end{aligned}$$

Ohm's law states that many substances, including many solids and some liquids, display this kind of behavior, at least for voltages that are not too large. The fact that Ohm's law is called a "law" should not be taken to mean that all materials obey it, or that it has the same fundamental importance as Newton's laws, for example. Materials are called *ohmic* or *nonohmic*, depending on whether they obey Ohm's law.

If objects of the same size and shape made from two different ohmic materials have different resistances, we can say that one material is more resistive than the other, or equivalently that it is less conductive. Materials, such as metals, that are very conductive are said to be good *conductors*. Those that are extremely poor conductors, for example wood or rubber, are classified as *insulators*. There is no sharp distinction between the two classes of materials. Some, such as silicon, lie midway between the two extremes, and are called *semiconductors*.

On an intuitive level, we can understand the idea of resistance by making the sounds "hhhhh" and "fffff." To make air flow out of your mouth, you use your diaphragm to compress the air in your



f / Four objects made of the same substance have different resistances.

chest. The pressure difference between your chest and the air outside your mouth is analogous to a voltage difference. When you make the “h” sound, you form your mouth and throat in a way that allows air to flow easily. The large flow of air is like a large current. Dividing by a large current in the definition of resistance means that we get a small resistance. We say that the small resistance of your mouth and throat allows a large current to flow. When you make the “f” sound, you increase the resistance and cause a smaller current to flow.

Note that although the resistance of an object depends on the substance it is made of, we cannot speak simply of the “resistance of gold” or the “resistance of wood.” Figure f shows four examples of objects that have had wires attached at the ends as electrical connections. If they were made of the same substance, they would all nevertheless have different resistances because of their different sizes and shapes. A more detailed discussion will be more natural in the context of the following chapter, but it should not be too surprising that the resistance of $f/2$ will be greater than that of $f/1$ — the image of water flowing through a pipe, however incorrect, gives us the right intuition. Object $f/3$ will have a smaller resistance than $f/1$ because the charged particles have less of it to get through.

Superconductors

All materials display some variation in resistance according to temperature (a fact that is used in thermostats to make a thermometer that can be easily interfaced to an electric circuit). More spectacularly, most metals have been found to exhibit a sudden change to *zero* resistance when cooled to a certain critical temperature. They are then said to be superconductors. Theoretically, superconductors should make a great many exciting devices possible, for example coiled-wire magnets that could be used to levitate trains. In practice, the critical temperatures of all metals are very low, and the resulting need for extreme refrigeration has made their use uneconomical except for such specialized applications as particle accelerators for physics research.

But scientists have recently made the surprising discovery that certain ceramics are superconductors at less extreme temperatures. The technological barrier is now in finding practical methods for making wire out of these brittle materials. Wall Street is currently investing billions of dollars in developing superconducting devices for cellular phone relay stations based on these materials. In 2001, the city of Copenhagen replaced a short section of its electrical power trunks with superconducting cables, and they are now in operation and supplying power to customers.

There is currently no satisfactory theory of superconductivity in general, although superconductivity in metals is understood fairly



g / A superconducting segment of the ATLAS accelerator at Argonne National Laboratory near Chicago. It is used to accelerate beams of ions to a few percent of the speed of light for nuclear physics research. The shiny silver-colored surfaces are made of the element niobium, which is a superconductor at relatively high temperatures compared to other metals — relatively high meaning the temperature of liquid helium! The beam of ions passes through the holes in the two small cylinders on the ends of the curved rods. Charge is shuffled back and forth between them at a frequency of 12 million cycles per second, so that they take turns being positive and negative. The positively charged beam consists of short spurts, each timed so that when it is in one of the segments it will be pulled forward by negative charge on the cylinder in front of it and pushed forward by the positively charged one behind. The huge currents involved would quickly melt any metal that was not superconducting, but in a superconductor they produce no heat at all.

well. Unfortunately I have yet to find a fundamental explanation of superconductivity in metals that works at the introductory level.

Finding charge given current

example 8

▷ In the segment of the ATLAS accelerator shown in figure g, the current flowing back and forth between the two cylinders is given by $I = a \cos bt$. What is the charge on one of the cylinders as a function of time? ▷ We are given the current and want to find the charge, i.e., we are given the derivative and we want to find the original function that would give that derivative. This means we need to integrate:

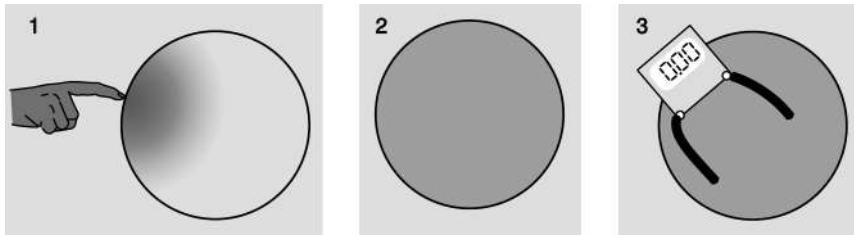
$$\begin{aligned} q &= \int I dt \\ &= \int a \cos bt dt \\ &= \frac{a}{b} \sin bt + q_0, \end{aligned}$$

where q_0 is a constant of integration.

We can interpret this in order to explain why a superconductor needs to be used. The constant b must be very large, since the current is supposed to oscillate back and forth millions of times a second. Looking at the final result, we see that if b is a very large number, and q is to be a significant amount of charge, then a must be a very large number as well. If a is numerically large, then the current must be very large, so it would heat the accelerator too much if it was flowing through an ordinary conductor.

Constant potential throughout a conductor

The idea of a superconductor leads us to the question of how we should expect an object to behave if it is made of a very good conductor. Superconductors are an extreme case, but often a metal wire can be thought of as a perfect conductor, for example if the parts of the circuit other than the wire are made of much less conductive materials. What happens if R equals zero in the equation $R = \Delta V/I$? The result of dividing two numbers can only be zero if the number on top equals zero. This tells us that if we pick any two points in a perfect conductor, the voltage difference between them must be zero. In other words, the entire conductor must be at the same potential.



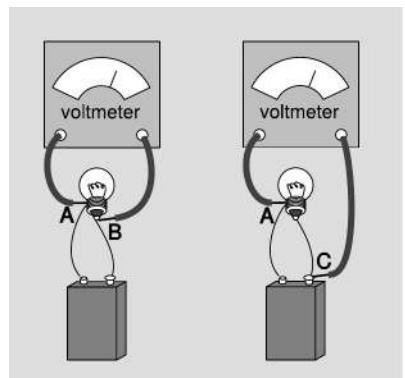
Constant potential means that no work would be done on a charge as it moved from one point in the conductor to another. If zero work was done only along a certain path between two specific points, it might mean that positive work was done along part of the path and negative work along the rest, resulting in a cancellation. But there is no way that the work could come out to be zero for all possible paths unless the electrical force on a charge was in fact zero at every point. Suppose, for example, that you build up a static charge by scuffing your feet on a carpet, and then you deposit some of that charge onto a doorknob, which is a good conductor. How can all that charge be in the doorknob without creating any electrical force at any point inside it? The only possible answer is that the charge moves around until it has spread itself into just the right configuration so that the forces exerted by all the little bits of excess surface charge on any charged particle within the doorknob exactly cancel out.

We can explain this behavior if we assume that the charge placed on the doorknob eventually settles down into a stable equilibrium. Since the doorknob is a conductor, the charge is free to move through it. If it was free to move and any part of it did experience a nonzero total force from the rest of the charge, then it would move, and we would not have an equilibrium.

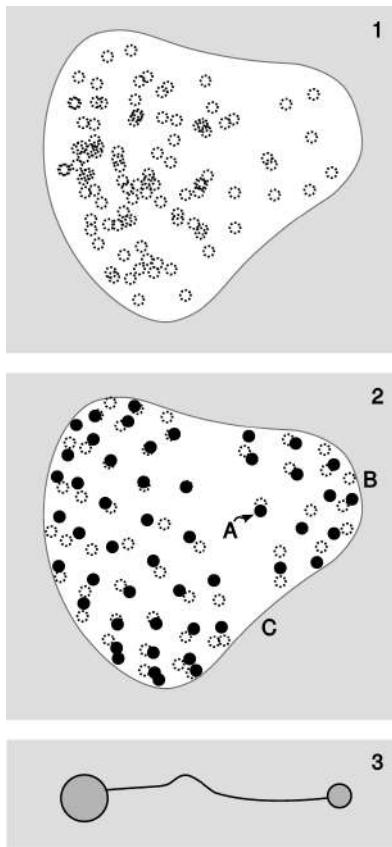
Excess charge placed on a conductor, once it reaches its equilibrium configuration, is entirely on the surface, not on the interior. This should be intuitively reasonable in figure h, for example, since the charges are all repelling each other. A proof is given in example 38 on p. 647.

Since wires are good conductors, constancy of potential throughout a conductor provides a convenient freedom in hooking up a voltmeter to a circuit. In figure i, points B and C are on the same piece of conducting wire, so $V_B = V_C$. Measuring $V_B - V_A$ gives the same result as measuring $V_C - V_A$.

- h / 1. The finger deposits charges on the solid, spherical, metal doorknob and is then withdrawn.
 2. Almost instantaneously, the charges' mutual repulsion makes them redistribute themselves uniformly on the surface of the sphere. The only excess charge is on the surface; charges do exist in the atoms that form the interior of the sphere, but they are balanced. Charges on the interior feel zero total electrical force from the ones at the surface. Charges at the surface experience a net outward repulsion, but this is canceled out by the force that keeps them from escaping into the air.
 3. A voltmeter shows zero difference in voltage between any two points on the interior or surface of the sphere. If the voltage difference wasn't zero, then energy could be released by the flow of charge from one point to the other; this only happens before equilibrium is reached.



i / The voltmeter doesn't care which of these setups you use.



j / Example 9. In 1 and 2, charges that are visible on the front surface of the conductor are shown as solid dots; the others would have to be seen through the conductor, which we imagine is semi-transparent.

The lightning rod

example 9

Suppose you have a pear-shaped conductor like the one in figure j/1. Since the pear is a conductor, there are free charges everywhere inside it. Panels 1 and 2 of the figure show a computer simulation with 100 identical electric charges. In 1, the charges are released at random positions inside the pear. Repulsion causes them all to fly outward onto the surface and then settle down into an orderly but nonuniform pattern.

We might not have been able to guess the pattern in advance, but we can verify that some of its features make sense. For example, charge A has more neighbors on the right than on the left, which would tend to make it accelerate off to the left. But when we look at the picture as a whole, it appears reasonable that this is prevented by the larger number of more distant charges on its left than on its right.

There also seems to be a pattern to the nonuniformity: the charges collect more densely in areas like B, where the surface is strongly curved, and less densely in flatter areas like C.

To understand the reason for this pattern, consider j/3. Two conducting spheres are connected by a conducting wire. Since the whole apparatus is conducting, it must all be at one potential. As shown in problem 37 on p. 572, the density of charge is greater on the smaller sphere. This is an example of a more general fact observed in j/2, which is that the charge on a conductor packs itself more densely in areas that are more sharply curved.

Similar reasoning shows why Benjamin Franklin used a sharp tip when he invented the lightning rod. The charged stormclouds induce positive and negative charges to move to opposite ends of the rod. At the pointed upper end of the rod, the charge tends to concentrate at the point, and this charge attracts the lightning. The same effect can sometimes be seen when a scrap of aluminum foil is inadvertently put in a microwave oven. Modern experiments (Moore *et al.*, *Journal of Applied Meteorology* 39 (1999) 593) show that although a sharp tip is best at starting a spark, a more moderate curve, like the right-hand tip of the pear in this example, is better at successfully sustaining the spark for long enough to connect a discharge to the clouds.

Short circuits

So far we have been assuming a perfect conductor. What if it is a good conductor, but not a perfect one? Then we can solve for $\Delta V = IR$. An ordinary-sized current will make a very small result when we multiply it by the resistance of a good conductor such as a metal wire. The potential throughout the wire will then be nearly constant. If, on the other hand, the current is extremely large, we can have a significant voltage difference. This is what happens in a

k / Short-circuiting a battery. Warning: you can burn yourself this way or start a fire! If you want to try this, try making the connection only very briefly, use a low-voltage battery, and avoid touching the battery or the wire, both of which will get hot.

short-circuit: a circuit in which a low-resistance pathway connects the two sides of a voltage source. Note that this is much more specific than the popular use of the term to indicate any electrical malfunction at all. If, for example, you short-circuit a 9-volt battery as shown in figure k, you will produce perhaps a thousand amperes of current, leading to a very large value of $P = I\Delta V$. The wire gets hot!

self-check A

What would happen to the battery in this kind of short circuit? ▶

Answer, p. 1062

At this stage, most students have a hard time understanding why resistors would be used inside a radio or a computer. We obviously want a lightbulb or an electric stove to have a circuit element that resists the flow of electricity and heats up, but heating is undesirable in radios and computers. Without going too far afield, let's use a mechanical analogy to get a general idea of why a resistor would be used in a radio.

The main parts of a radio receiver are an antenna, a tuner for selecting the frequency, and an amplifier to strengthen the signal sufficiently to drive a speaker. The tuner resonates at the selected frequency, just as in the examples of mechanical resonance discussed in 3. The behavior of a mechanical resonator depends on three things: its inertia, its stiffness, and the amount of friction or damping. The first two parameters locate the peak of the resonance curve, while the damping determines the width of the resonance. In the radio tuner we have an electrically vibrating system that resonates at a particular frequency. Instead of a physical object moving back and forth, these vibrations consist of electrical currents that flow first in one direction and then in the other. In a mechanical system, damping means taking energy out of the vibration in the form of heat, and exactly the same idea applies to an electrical system: the resistor supplies the damping, and therefore controls the width of the resonance. If we set out to eliminate all resistance in the tuner circuit, by not building in a resistor and by somehow getting rid of all the inherent electrical resistance of the wires, we would have a useless radio. The tuner's resonance would be so narrow that we could never get close enough to the right frequency to bring in the station. The roles of inertia and stiffness are played by other circuit elements we have not discussed (a capacitor and a coil).

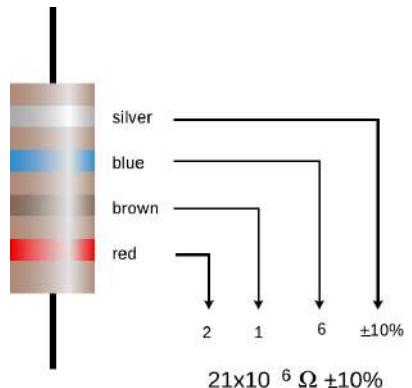
Resistors

Inside any electronic gadget you will see quite a few little circuit elements like the one shown below. These *resistors* are simply a cylinder of ohmic material with wires attached to the end.

Many electrical devices are based on electrical resistance and Ohm's law, even if they do not have little components in them that

black	0
brown	1
red	2
orange	3
yellow	4
green	5
blue	6
violet	7
gray	8
white	9
silver	$\pm 10\%$
gold	$\pm 5\%$

l / Color codes used on resistors.



m / An example of a resistor with a color code.



n / The symbol used in schematics to represent a resistor.

look like the usual resistor. The following are some examples.

Lightbulb

There is nothing special about a lightbulb filament — you can easily make a lightbulb by cutting a narrow waist into a metallic gum wrapper and connecting the wrapper across the terminals of a 9-volt battery. The trouble is that it will instantly burn out. Edison solved this technical challenge by encasing the filament in an evacuated bulb, which prevented burning, since burning requires oxygen.

Polygraph

The polygraph, or “lie detector,” is really just a set of meters for recording physical measures of the subject’s psychological stress, such as sweating and quickened heartbeat. The real-time sweat measurement works on the principle that dry skin is a good insulator, but sweaty skin is a conductor. Of course a truthful subject may become nervous simply because of the situation, and a practiced liar may not even break a sweat. The method’s practitioners claim that they can tell the difference, but you should think twice before allowing yourself to be polygraph tested. Most U.S. courts exclude all polygraph evidence, but some employers attempt to screen out dishonest employees by polygraph testing job applicants, an abuse that ranks with such pseudoscience as handwriting analysis.

Fuse

A fuse is a device inserted in a circuit tollbooth-style in the same manner as an ammeter. It is simply a piece of wire made of metals having a relatively low melting point. If too much current passes through the fuse, it melts, opening the circuit. The purpose is to make sure that the building’s wires do not carry so much current that they themselves will get hot enough to start a fire. Most modern houses use circuit breakers instead of fuses, although fuses are still common in cars and small devices. A circuit breaker is a switch operated by a coiled-wire magnet, which opens the circuit when enough current flows. The advantage is that once you turn off some of the appliances that were sucking up too much current, you can immediately flip the switch closed. In the days of fuses, one might get caught without a replacement fuse, or even be tempted to stuff aluminum foil in as a replacement, defeating the safety feature.

Voltmeter

A voltmeter is nothing more than an ammeter with an additional high-value resistor through which the current is also forced to flow. Ohm’s law states that the current through the resistor is related directly to the voltage difference across it, so the meter can

be calibrated in units of volts based on the known value of the resistor. The voltmeter's two probes are touched to the two locations in a circuit between which we wish to measure the voltage difference, o/2. Note how cumbersome this type of drawing is, and how difficult it can be to tell what is connected to what. This is why electrical drawing are usually shown in schematic form. Figure o/3 is a schematic representation of figure o/2.

The setups for measuring current and voltage are different. When we are measuring current, we are finding “how much stuff goes through,” so we place the ammeter where all the current is forced to go through it. Voltage, however, is not “stuff that goes through,” it is a measure of electrical energy. If an ammeter is like the meter that measures your water use, a voltmeter is like a measuring stick that tells you how high a waterfall is, so that you can determine how much energy will be released by each kilogram of falling water. We do not want to force the water to go through the measuring stick! The arrangement in figure o/3 is a *parallel* circuit: one in there are “forks in the road” where some of the current will flow one way and some will flow the other. Figure o/4 is said to be wired in *series*: all the current will visit all the circuit elements one after the other. We will deal with series and parallel circuits in more detail in the following chapter.

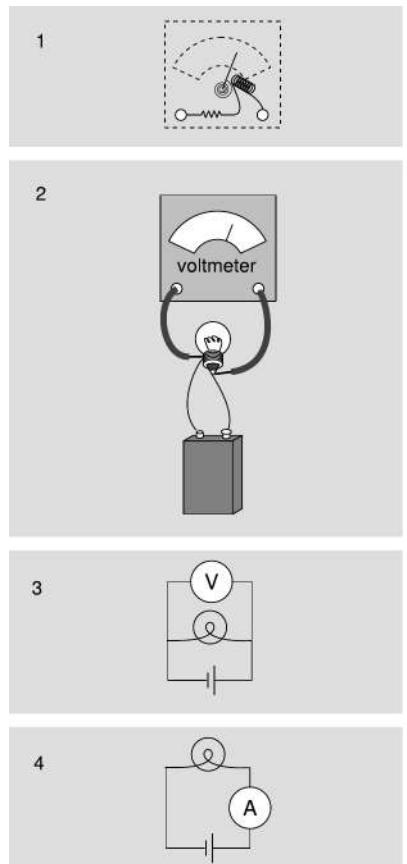
If you inserted a voltmeter incorrectly, in series with the bulb and battery, its large internal resistance would cut the current down so low that the bulb would go out. You would have severely disturbed the behavior of the circuit by trying to measure something about it.

Incorrectly placing an ammeter in parallel is likely to be even more disconcerting. The ammeter has nothing but wire inside it to provide resistance, so given the choice, most of the current will flow through it rather than through the bulb. So much current will flow through the ammeter, in fact, that there is a danger of burning out the battery or the meter or both! For this reason, most ammeters have fuses or circuit breakers inside. Some models will trip their circuit breakers and make an audible alarm in this situation, while others will simply blow a fuse and stop working until you replace it.

Discussion Questions

A In figure o/1, would it make any difference in the voltage measurement if we touched the voltmeter's probes to different points along the same segments of wire?

B Explain why it would be incorrect to define resistance as the amount of charge the resistor allows to flow.



o / 1. A simplified diagram of how a voltmeter works. 2. Measuring the voltage difference across a lightbulb. 3. The same setup drawn in schematic form. 4. The setup for measuring current is different.

9.1.5 Current-conducting properties of materials

Ohm's law has a remarkable property, which is that current will flow even in response to a voltage difference that is as small as we care to make it. In the analogy of pushing a crate across a floor, it is as though even a flea could slide the crate across the floor, albeit at some very low speed. The flea cannot do this because of static friction, which we can think of as an effect arising from the tendency of the microscopic bumps and valleys in the crate and floor to lock together. The fact that Ohm's law holds for nearly all solids has an interesting interpretation: at least some of the electrons are not "locked down" at all to any specific atom.

More generally we can ask how charge actually flows in various solids, liquids, and gases. This will lead us to the explanations of many interesting phenomena, including lightning, the bluish crust that builds up on the terminals of car batteries, and the need for electrolytes in sports drinks.

Solids

In atomic terms, the defining characteristic of a solid is that its atoms are packed together, and the nuclei cannot move very far from their equilibrium positions. It makes sense, then, that electrons, not ions, would be the charge carriers when currents flow in solids. This fact was established experimentally by Tolman and Stewart, in an experiment in which they spun a large coil of wire and then abruptly stopped it. They observed a current in the wire immediately after the coil was stopped, which indicated that charged particles that were not permanently locked to a specific atom had continued to move because of their own inertia, even after the material of the wire in general stopped. The direction of the current showed that it was negatively charged particles that kept moving. The current only lasted for an instant, however; as the negatively charged particles collected at the downstream end of the wire, farther particles were prevented joining them due to their electrical repulsion, as well as the attraction from the upstream end, which was left with a net positive charge. Tolman and Stewart were even able to determine the mass-to-charge ratio of the particles. We need not go into the details of the analysis here, but particles with high mass would be difficult to decelerate, leading to a stronger and longer pulse of current, while particles with high charge would feel stronger electrical forces decelerating them, which would cause a weaker and shorter pulse. The mass-to-charge ratio thus determined was consistent with the m/q of the electron to within the accuracy of the experiment, which essentially established that the particles were electrons.

The fact that only electrons carry current in solids, not ions, has many important implications. For one thing, it explains why wires don't fray or turn to dust after carrying current for a long time. Electrons are very small (perhaps even pointlike), and it is easy to

imagine them passing between the cracks among the atoms without creating holes or fractures in the atomic framework. For those who know a little chemistry, it also explains why all the best conductors are on the left side of the periodic table. The elements in that area are the ones that have only a very loose hold on their outermost electrons.

Gases

The molecules in a gas spend most of their time separated from each other by significant distances, so it is not possible for them to conduct electricity the way solids do, by handing off electrons from atom to atom. It is therefore not surprising that gases are good insulators.

Gases are also usually nonohmic. As opposite charges build up on a stormcloud and the ground below, the voltage difference becomes greater and greater. Zero current flows, however, until finally the voltage reaches a certain threshold and we have an impressive example of what is known as a spark or electrical discharge. If air was ohmic, the current between the cloud and the ground would simply increase steadily as the voltage difference increased, rather than being zero until a threshold was reached. This behavior can be explained as follows. At some point, the electrical forces on the air electrons and nuclei of the air molecules become so strong that electrons are ripped right off of some of the molecules. The electrons then accelerate toward either the cloud or the ground, whichever is positively charged, and the positive ions accelerate the opposite way. As these charge carriers accelerate, they strike and ionize other molecules, which produces a rapidly growing cascade.

Liquids

Molecules in a liquid are able to slide past each other, so ions as well as electrons can carry currents. Pure water is a poor conductor because the water molecules tend to hold onto their electrons strongly, and there are therefore not many electrons or ions available to move. Water can become quite a good conductor, however, with the addition of even a small amount of certain substances called electrolytes, which are typically salts. For example, if we add table salt, NaCl, to water, the NaCl molecules dissolve into Na^+ and Cl^- ions, which can then move and create currents. This is why electric currents can flow among the cells in our bodies: cellular fluid is quite salty. When we sweat, we lose not just water but electrolytes, so dehydration plays havoc with our cells' electrical systems. It is for this reason that electrolytes are included in sports drinks and formulas for rehydrating infants who have diarrhea.

Since current flow in liquids involves entire ions, it is not surprising that we can see physical evidence when it has occurred. For example, after a car battery has been in use for a while, the H_2SO_4

battery acid becomes depleted of hydrogen ions, which are the main charge carriers that complete the circuit on the inside of the battery. The leftover SO_4 then forms a visible blue crust on the battery posts.

Speed of currents and electrical signals

When I talk on the phone to my mother in law two thousand miles away, I do not notice any delay while the signal makes its way back and forth. Electrical signals therefore must travel very quickly, but how fast exactly? The answer is rather subtle. For the sake of concreteness, let's restrict ourselves to currents in metals, which consist of electrons.

The electrons themselves are only moving at speeds of perhaps a few thousand miles per hour, and their motion is mostly random thermal motion. This shows that the electrons in my phone cannot possibly be zipping back and forth between California and New York fast enough to carry the signals. Even if their thousand-mile-an-hour motion was organized rather than random, it would still take them many minutes to get there. Realistically, it will take the average electron even longer than that to make the trip. The current in the wire consists only of a slow overall drift, at a speed on the order of a few centimeters per second, superimposed on the more rapid random motion. We can compare this with the slow westward drift in the population of the U.S. If we could make a movie of the motion of all the people in the U.S. from outer space, and could watch it at high speed so that the people appeared to be scurrying around like ants, we would think that the motion was fairly random, and we would not immediately notice the westward drift. Only after many years would we realize that the number of people heading west over the Sierras had exceeded the number going east, so that California increased its share of the country's population.

So why are electrical signals so fast if the average drift speed of electrons is so slow? The answer is that a disturbance in an electrical system can move much more quickly than the charges themselves. It is as though we filled a pipe with golf balls and then inserted an extra ball at one end, causing a ball to fall out at the other end. The force propagated to the other end in a fraction of a second, but the balls themselves only traveled a few centimeters in that time.

Because the reality of current conduction is so complex, we often describe things using mental shortcuts that are technically incorrect. This is OK as long as we know that they are just shortcuts. For example, suppose the presidents of France and Russia shake hands, and the French politician has inadvertently picked up a positive electrical charge, which shocks the Russian. We may say that the excess positively charged particles in the French leader's body, which all repel each other, take the handshake as an opportunity to get farther apart by spreading out into two bodies rather than one. In

reality, it would be a matter of minutes before the ions in one person's body could actually drift deep into the other's. What really happens is that throughout the body of the recipient of the shock there are already various positive and negative ions which are free to move. Even before the perpetrator's charged hand touches the victim's sweaty palm, the charges in the shocker's body begin to repel the positive ions and attract the negative ions in the other person. The split-second sensation of shock is caused by the sudden jumping of the victim's ions by distances of perhaps a micrometer, this effect occurring simultaneously throughout the whole body, although more violently in the hand and arm, which are closer to the other person.

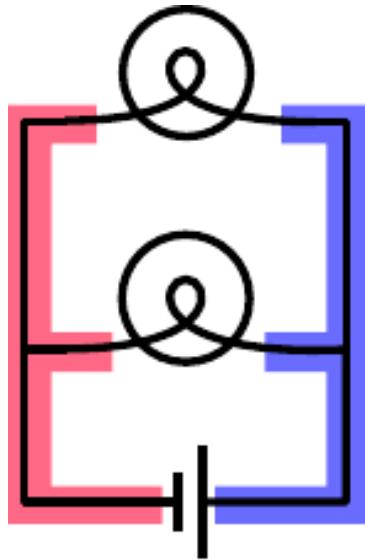
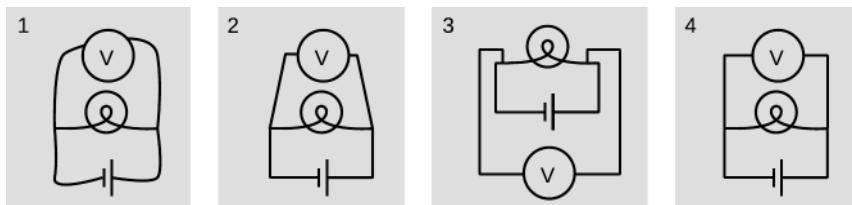
9.2 Parallel and series circuits

In section 9.1, we limited ourselves to relatively simple circuits, essentially nothing more than a battery and a single lightbulb. The purpose of this chapter is to introduce you to more complex circuits, containing multiple resistors or voltage sources in series, in parallel, or both.

9.2.1 Schematics

I see a chess position; Kasparov sees an interesting Ruy Lopez variation. To the uninitiated a schematic may look as unintelligible as Mayan hieroglyphs, but even a little bit of eye training can go a long way toward making its meaning leap off the page. A schematic is a stylized and simplified drawing of a circuit. The purpose is to eliminate as many irrelevant features as possible, so that the relevant ones are easier to pick out.

- a / 1. Wrong: The shapes of the wires are irrelevant. 2. Wrong: Right angles should be used. 3. Wrong: A simple pattern is made to look unfamiliar and complicated. 4. Right.



- b / The two shaded areas shaped like the letter "E" are both regions of constant voltage.

An example of an irrelevant feature is the physical shape, length, and diameter of a wire. In nearly all circuits, it is a good approximation to assume that the wires are perfect conductors, so that any piece of wire uninterrupted by other components has constant voltage throughout it. Changing the length of the wire, for instance, does not change this fact. (Of course if we used miles and miles of wire, as in a telephone line, the wire's resistance would start to add up, and its length would start to matter.) The shapes of the wires are likewise irrelevant, so we draw them with standardized, stylized shapes made only of vertical and horizontal lines with right-angle bends in them. This has the effect of making similar circuits look more alike and helping us to recognize familiar patterns, just as words in a newspaper are easier to recognize than handwritten ones. Figure a shows some examples of these concepts.

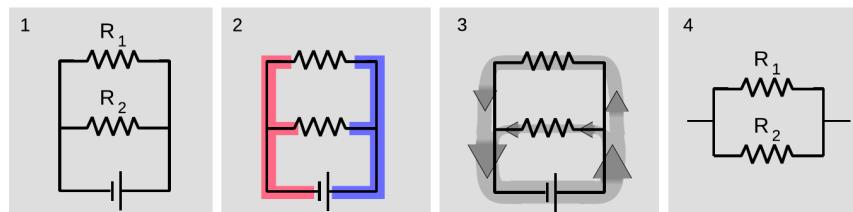
The most important first step in learning to read schematics is to learn to recognize contiguous pieces of wire which must have constant voltage throughout. In figure b, for example, the two shaded E-shaped pieces of wire must each have constant voltage. This focuses our attention on two of the main unknowns we'd like to be able to predict: the voltage of the left-hand E and the voltage of the one on the right.

9.2.2 Parallel resistances and the junction rule

One of the simplest examples to analyze is the parallel resistance circuit, of which figure b was an example. In general we may have unequal resistances R_1 and R_2 , as in c/1. Since there are only two constant-voltage areas in the circuit, c/2, all three components have the same voltage difference across them. A battery normally succeeds in maintaining the voltage differences across itself for which it was designed, so the voltage drops ΔV_1 and ΔV_2 across the resistors must both equal the voltage of the battery:

$$\Delta V_1 = \Delta V_2 = \Delta V_{battery}.$$

Each resistance thus feels the same voltage difference as if it was the only one in the circuit, and Ohm's law tells us that the amount of current flowing through each one is also the same as it would have been in a one-resistor circuit. This is why household electrical circuits are wired in parallel. We want every appliance to work the same, regardless of whether other appliances are plugged in or unplugged, turned on or switched off. (The electric company doesn't use batteries of course, but our analysis would be the same for any device that maintains a constant voltage.)



- c / 1. Two resistors in parallel.
- 2. There are two constant-voltage areas.
- 3. The current that comes out of the battery splits between the two resistors, and later reunites.
- 4. The two resistors in parallel can be treated as a single resistor with a smaller resistance value.

Of course the electric company can tell when we turn on every light in the house. How do they know? The answer is that we draw more current. Each resistance draws a certain amount of current, and the amount that has to be supplied is the sum of the two individual currents. The current is like a river that splits in half, c/3, and then reunites. The total current is

$$I_{total} = I_1 + I_2.$$

This is an example of a general fact called the junction rule:

In any circuit that is not storing or releasing charge, conservation of charge implies that the total current flowing out of any junction must be the same as the total flowing in.

Coming back to the analysis of our circuit, we apply Ohm's law

to each resistance, resulting in

$$\begin{aligned}I_{total} &= \Delta V/R_1 + \Delta V/R_2 \\&= \Delta V \left(\frac{1}{R_1} + \frac{1}{R_2} \right).\end{aligned}$$

As far as the electric company is concerned, your whole house is just one resistor with some resistance R , called the *equivalent resistance*. They would write Ohm's law as

$$I_{total} = \Delta V/R,$$

from which we can determine the equivalent resistance by comparison with the previous expression:

$$\begin{aligned}\frac{1}{R} &= \frac{1}{R_1} + \frac{1}{R_2} \\R &= \left(\frac{1}{R_1} + \frac{1}{R_2} \right)^{-1}\end{aligned}$$

[equivalent resistance of two resistors in parallel]

Two resistors in parallel, $c/4$, are equivalent to a single resistor with a value given by the above equation.

Two lamps on the same household circuit *example 10*

▷ You turn on two lamps that are on the same household circuit. Each one has a resistance of 1 ohm. What is the equivalent resistance, and how does the power dissipation compare with the case of a single lamp?

▷ The equivalent resistance of the two lamps in parallel is

$$\begin{aligned}R &= \left(\frac{1}{R_1} + \frac{1}{R_2} \right)^{-1} \\&= \left(\frac{1}{1\Omega} + \frac{1}{1\Omega} \right)^{-1} \\&= \left(1\Omega^{-1} + 1\Omega^{-1} \right)^{-1} \\&= \left(2\Omega^{-1} \right)^{-1} \\&= 0.5\Omega\end{aligned}$$

The voltage difference across the whole circuit is always the 110 V set by the electric company (it's alternating current, but that's irrelevant). The resistance of the whole circuit has been cut in half by turning on the second lamp, so a fixed amount of voltage will produce twice as much current. Twice the current flowing across the same voltage difference means twice as much power dissipation, which makes sense.

The cutting in half of the resistance surprises many students, since we are “adding more resistance” to the circuit by putting in the second lamp. Why does the equivalent resistance come out to be less than the resistance of a single lamp? This is a case where purely verbal reasoning can be misleading. A resistive circuit element, such as the filament of a lightbulb, is neither a perfect insulator nor a perfect conductor. Instead of analyzing this type of circuit in terms of “resistors,” i.e., partial insulators, we could have spoken of “conductors.” This example would then seem reasonable, since we “added more conductance,” but one would then have the incorrect expectation about the case of resistors in series, discussed in the following section.

Perhaps a more productive way of thinking about it is to use mechanical intuition. By analogy, your nostrils resist the flow of air through them, but having two nostrils makes it twice as easy to breathe.

Three resistors in parallel

example 11

- ▷ What happens if we have three or more resistors in parallel?
- ▷ This is an important example, because the solution involves an important technique for understanding circuits: breaking them down into smaller parts and then simplifying those parts. In the circuit d/1, with three resistors in parallel, we can think of two of the resistors as forming a single big resistor, d/2, with equivalent resistance

$$R_{12} = \left(\frac{1}{R_1} + \frac{1}{R_2} \right)^{-1}.$$

We can then simplify the circuit as shown in d/3, so that it contains only two resistances. The equivalent resistance of the whole circuit is then given by

$$R_{123} = \left(\frac{1}{R_{12}} + \frac{1}{R_3} \right)^{-1}.$$

Substituting for R_{12} and simplifying, we find the result

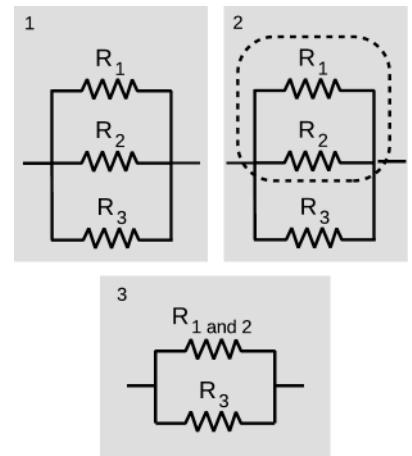
$$R_{123} = \left(\frac{1}{R_1} + \frac{1}{R_2} + \frac{1}{R_3} \right)^{-1},$$

which you probably could have guessed. The interesting point here is the divide-and-conquer concept, not the mathematical result.

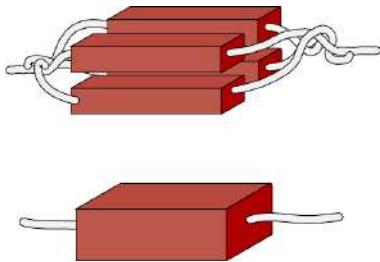
An arbitrary number of identical resistors in parallel example 12

- ▷ What is the resistance of N identical resistors in parallel?
- ▷ Generalizing the results for two and three resistors, we have

$$R_N = \left(\frac{1}{R_1} + \frac{1}{R_2} + \dots \right)^{-1},$$



d / Three resistors in parallel.



e / Uniting four resistors in parallel is equivalent to making a single resistor with the same length but four times the cross-sectional area. The result is to make a resistor with one quarter the resistance.

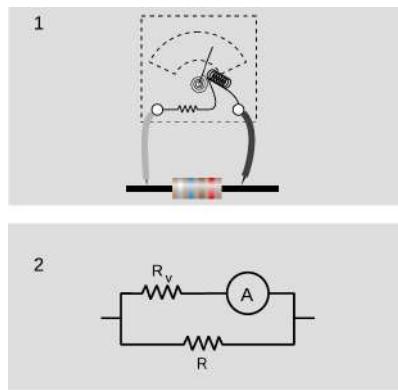
where “...” means that the sum includes all the resistors. If all the resistors are identical, this becomes

$$R_N = \left(\frac{N}{R} \right)^{-1}$$

$$= \frac{R}{N}$$

Dependence of resistance on cross-sectional area example 13
We have alluded briefly to the fact that an object's electrical resistance depends on its size and shape, but now we are ready to begin making more mathematical statements about it. As suggested by figure e, increasing a resistor's cross-sectional area is equivalent to adding more resistors in parallel, which will lead to an overall decrease in resistance. Any real resistor with straight, parallel sides can be sliced up into a large number of pieces, each with cross-sectional area of, say, $1 \mu\text{m}^2$. The number, N , of such slices is proportional to the total cross-sectional area of the resistor, and by application of the result of the previous example we therefore find that the resistance of an object is inversely proportional to its cross-sectional area.

f / A fat pipe has less resistance than a skinny pipe.



g / A voltmeter is really an ammeter with an internal resistor. When we measure the voltage difference across a resistor, 1, we are really constructing a parallel resistance circuit, 2.

An analogous relationship holds for water pipes, which is why high-flow trunk lines have to have large cross-sectional areas. To make lots of water (current) flow through a skinny pipe, we'd need an impractically large pressure (voltage) difference.

Incorrect readings from a voltmeter example 14
A voltmeter is really just an ammeter with an internal resistor, and we use a voltmeter in parallel with the thing that we're trying to measure the voltage difference across. This means that any time we measure the voltage drop across a resistor, we're essentially putting two resistors in parallel. The ammeter inside the voltmeter can be ignored for the purpose of analyzing what how current flows in the circuit, since it is essentially just some coiled-up wire with a very low resistance.

Now if we are carrying out this measurement on a resistor that is part of a larger circuit, we have changed the behavior of the circuit through our act of measuring. It is as though we had modified the circuit by replacing the resistance R with the smaller equivalent resistance of R and R_v in parallel. It is for this reason that voltmeters are built with the largest possible internal resistance. As a numerical example, if we use a voltmeter with an internal resistance of $1\text{ M}\Omega$ to measure the voltage drop across a one-ohm resistor, the equivalent resistance is $0.999999\text{ }\Omega$, which is not different enough to make any difference. But if we tried to use the same voltmeter to measure the voltage drop across a $2 - \text{M}\Omega$ resistor, we would be reducing the resistance of that part of the circuit by a factor of three, which would produce a drastic change in the behavior of the whole circuit.

This is the reason why you can't use a voltmeter to measure the voltage difference between two different points in mid-air, or between the ends of a piece of wood. This is by no means a stupid thing to want to do, since the world around us is not a constant-voltage environment, the most extreme example being when an electrical storm is brewing. But it will not work with an ordinary voltmeter because the resistance of the air or the wood is many gigaohms. The effect of waving a pair of voltmeter probes around in the air is that we provide a reuniting path for the positive and negative charges that have been separated — through the voltmeter itself, which is a good conductor compared to the air. This reduces to zero the voltage difference we were trying to measure.

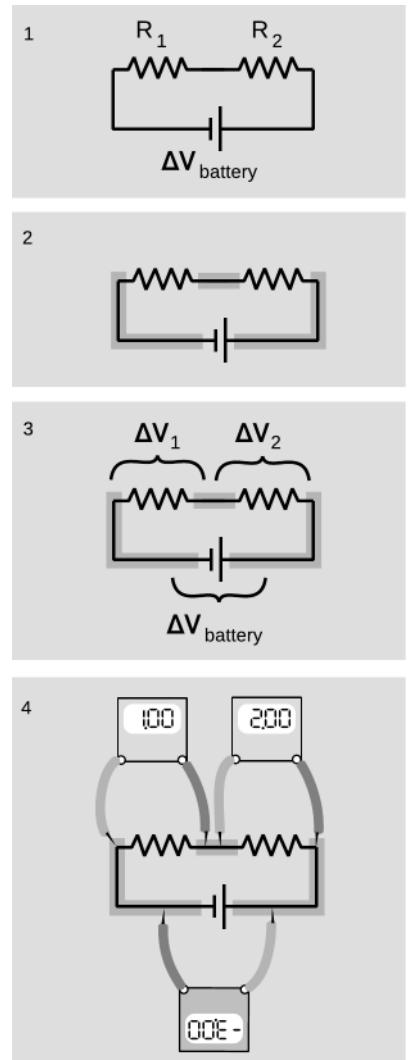
In general, a voltmeter that has been set up with an open circuit (or a very large resistance) between its probes is said to be “floating.” An old-fashioned analog voltmeter of the type described here will read zero when left floating, the same as when it was sitting on the shelf. A floating digital voltmeter usually shows an error message.

9.2.3 Series resistances

The two basic circuit layouts are parallel and series, so a pair of resistors in series, h/1, is another of the most basic circuits we can make. By conservation of charge, all the current that flows through one resistor must also flow through the other (as well as through the battery):

$$I_1 = I_2.$$

The only way the information about the two resistance values is going to be useful is if we can apply Ohm's law, which will relate the resistance of each resistor to the current flowing through it and the voltage difference across it. Figure h/2 shows the three constant-voltage areas. Voltage differences are more physically significant than voltages, so we define symbols for the voltage differences across



h / 1. A battery drives current through two resistors in series. 2. There are three constant-voltage regions. 3. The three voltage differences are related. 4. If the meter crab-walks around the circuit without flipping over or crossing its legs, the resulting voltages have plus and minus signs that make them add up to zero.

the two resistors in figure h/3.

We have three constant-voltage areas, with symbols for the difference in voltage between every possible pair of them. These three voltage differences must be related to each other. It is as though I tell you that Fred is a foot taller than Ginger, Ginger is a foot taller than Sally, and Fred is two feet taller than Sally. The information is redundant, and you really only needed two of the three pieces of data to infer the third. In the case of our voltage differences, we have

$$|\Delta V_1| + |\Delta V_2| = |\Delta V_{battery}|.$$

The absolute value signs are because of the ambiguity in how we define our voltage differences. If we reversed the two probes of the voltmeter, we would get a result with the opposite sign. Digital voltmeters will actually provide a minus sign on the screen if the wire connected to the “V” plug is lower in voltage than the one connected to the “COM” plug. Analog voltmeters pin the needle against a peg if you try to use them to measure negative voltages, so you have to fiddle to get the leads connected the right way, and then supply any necessary minus sign yourself.

Figure h/4 shows a standard way of taking care of the ambiguity in signs. For each of the three voltage measurements around the loop, we keep the same probe (the darker one) on the clockwise side. It is as though the voltmeter was sidling around the circuit like a crab, without ever “crossing its legs.” With this convention, the relationship among the voltage drops becomes

$$\Delta V_1 + \Delta V_2 = -\Delta V_{battery},$$

or, in more symmetrical form,

$$\Delta V_1 + \Delta V_2 + \Delta V_{battery} = 0.$$

More generally, this is known as the loop rule for analyzing circuits:

Assuming the standard convention for plus and minus signs, the sum of the voltage drops around any closed loop in a DC circuit must be zero.

Looking for an exception to the loop rule would be like asking for a hike that would be downhill all the way and that would come back to its starting point!

For the circuit we set out to analyze, the equation

$$\Delta V_1 + \Delta V_2 + \Delta V_{battery} = 0$$

can now be rewritten by applying Ohm's law to each resistor:

$$I_1 R_1 + I_2 R_2 + \Delta V_{battery} = 0.$$

The currents are the same, so we can factor them out:

$$I (R_1 + R_2) + \Delta V_{battery} = 0,$$

and this is the same result we would have gotten if we had been analyzing a one-resistor circuit with resistance $R_1 + R_2$. Thus the equivalent resistance of resistors in series equals the sum of their resistances.

Two lightbulbs in series

example 15

- ▷ If two identical lightbulbs are placed in series, how do their brightnesses compare with the brightness of a single bulb?
- ▷ Taken as a whole, the pair of bulbs act like a doubled resistance, so they will draw half as much current from the wall. Each bulb will be dimmer than a single bulb would have been.

The total power dissipated by the circuit is $I\Delta V$. The voltage drop across the whole circuit is the same as before, but the current is halved, so the two-bulb circuit draws half as much total power as the one-bulb circuit. Each bulb draws one-quarter of the normal power.

Roughly speaking, we might expect this to result in one quarter the light being produced by each bulb, but in reality lightbulbs waste quite a high percentage of their power in the form of heat and wavelengths of light that are not visible (infrared and ultraviolet). Less light will be produced, but it's hard to predict exactly how much less, since the efficiency of the bulbs will be changed by operating them under different conditions.

More than two equal resistances in series

example 16

By straightforward application of the divide-and-conquer technique discussed in the previous section, we find that the equivalent resistance of N identical resistances R in series will be NR .

Dependence of resistance on length

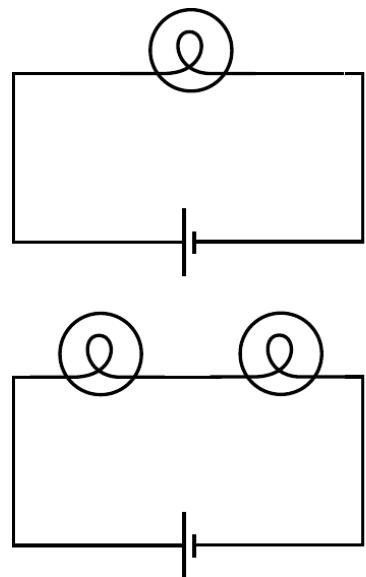
example 17

In the previous section, we proved that resistance is inversely proportional to cross-sectional area. By equivalent reason about resistances in series, we find that resistance is proportional to length. Analogously, it is harder to blow through a long straw than through a short one.

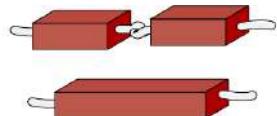
Putting the two arguments together, we find that the resistance of an object with straight, parallel sides is given by

$$R = (\text{constant}) \cdot L/A$$

The proportionality constant is called the resistivity, and it depends only on the substance of which the object is made. A resistivity measurement could be used, for instance, to help identify a sample of an unknown substance.



i / Example 15.



j / Doubling the length of a resistor is like putting two resistors in series. The resistance is doubled.

Choice of high voltage for power lines

example 18

Thomas Edison got involved in a famous technological controversy over the voltage difference that should be used for electrical power lines. At this time, the public was unfamiliar with electricity, and easily scared by it. The president of the United States, for instance, refused to have electrical lighting in the White House when it first became commercially available because he considered it unsafe, preferring the known fire hazard of oil lamps to the mysterious dangers of electricity. Mainly as a way to overcome public fear, Edison believed that power should be transmitted using small voltages, and he publicized his opinion by giving demonstrations at which a dog was lured into position to be killed by a large voltage difference between two sheets of metal on the ground. (Edison's opponents also advocated alternating current rather than direct current, and AC is more dangerous than DC as well. As we will discuss later, AC can be easily stepped up and down to the desired voltage level using a device called a transformer.)

Now if we want to deliver a certain amount of power P_L to a load such as an electric lightbulb, we are constrained only by the equation $P_L = I\Delta V_L$. We can deliver any amount of power we wish, even with a low voltage, if we are willing to use large currents. Modern electrical distribution networks, however, use dangerously high voltage differences of tens of thousands of volts. Why did Edison lose the debate?

It boils down to money. The electric company must deliver the amount of power P_L desired by the customer through a transmission line whose resistance R_T is fixed by economics and geography. The same current flows through both the load and the transmission line, dissipating power usefully in the former and wastefully in the latter. The efficiency of the system is

$$\begin{aligned}\text{efficiency} &= \frac{\text{power paid for by the customer}}{\text{power paid for by the utility}} \\ &= \frac{P_L}{P_L + P_T} \\ &= \frac{1}{1 + P_T/P_L}\end{aligned}$$

Putting ourselves in the shoes of the electric company, we wish to get rid of the variable P_T , since it is something we control only indirectly by our choice of ΔV_T and I . Substituting $P_T = I\Delta V_T$, we find

$$\text{efficiency} = \frac{1}{1 + \frac{I\Delta V_T}{P_L}}$$

We assume the transmission line (but not necessarily the load) is

ohmic, so substituting $\Delta V_T = IR_T$ gives

$$\text{efficiency} = \frac{1}{1 + \frac{I^2 R_T}{P_L}}$$

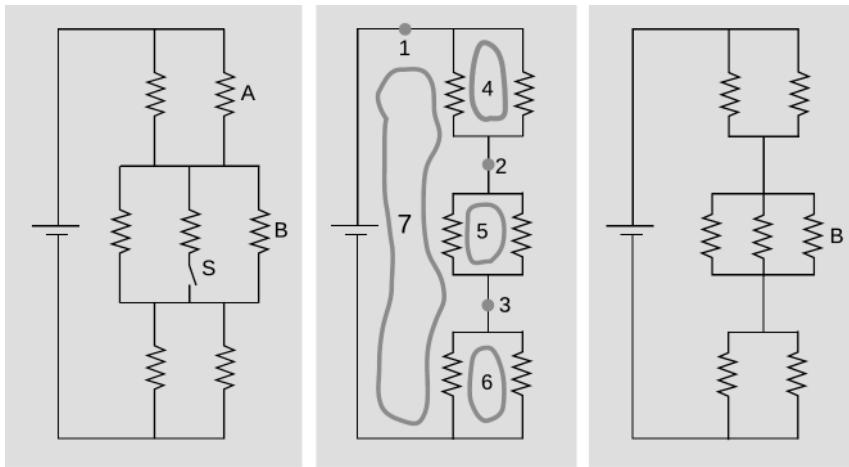
This quantity can clearly be maximized by making I as small as possible, since we will then be dividing by the smallest possible quantity on the bottom of the fraction. A low-current circuit can only deliver significant amounts of power if it uses high voltages, which is why electrical transmission systems use dangerous high voltages.

Getting killed by your ammeter

example 19

As with a voltmeter, an ammeter can give erroneous readings if it is used in such a way that it changes the behavior of the circuit. An ammeter is used in series, so if it is used to measure the current through a resistor, the resistor's value will effectively be changed to $R + R_a$, where R_a is the resistance of the ammeter. Ammeters are designed with very low resistances in order to make it unlikely that $R + R_a$ will be significantly different from R .

In fact, the real hazard is death, not a wrong reading! Virtually the only circuits whose resistances are significantly less than that of an ammeter are those designed to carry huge currents. An ammeter inserted in such a circuit can easily melt. When I was working at a laboratory funded by the Department of Energy, we got periodic bulletins from the DOE safety office about serious accidents at other sites, and they held a certain ghoulish fascination. One of these was about a DOE worker who was completely incinerated by the explosion created when he inserted an ordinary Radio Shack ammeter into a high-current circuit. Later estimates showed that the heat was probably so intense that the explosion was a ball of plasma — a gas so hot that its atoms have been ionized.



k / Example 20.

A complicated circuit

example 20

▷ All seven resistors in the left-hand panel of figure k are identical. Initially, the switch S is open as shown in the figure, and the current through resistor A is I_0 . The switch is then closed. Find the current through resistor B, after the switch is closed, in terms of I_0 .

▷ The second panel shows the circuit redrawn for simplicity, in the initial condition with the switch open. When the switch is open, no current can flow through the central resistor, so we may as well ignore it. I've also redrawn the junctions, without changing what's connected to what. This is the kind of mental rearranging that you'll eventually learn to do automatically from experience with analyzing circuits. The redrawn version makes it easier to see what's happening with the current. Charge is conserved, so any charge that flows past point 1 in the circuit must also flow past points 2 and 3. This would have been harder to reason about by applying the junction rule to the original version, which appears to have nine separate junctions.

In the new version, it's also clear that the circuit has a great deal of symmetry. We could flip over each parallel pair of identical resistors without changing what's connected to what, so that makes it clear that the voltage drops and currents must be equal for the members of each pair. We can also prove this by using the loop rule. The loop rule says that the two voltage drops in loop 4 must be equal, and similarly for loops 5 and 6. Since the resistors obey Ohm's law, equal voltage drops across them also imply equal currents. That means that when the current at point 1 comes to the top junction, exactly half of it goes through each resistor. Then the current reunites at 2, splits between the next pair, and so on. We conclude that each of the six resistors in the circuit experiences the same voltage drop and the same current. Applying the loop rule to loop 7, we find that the sum of the three voltage drops across the three left-hand resistors equals the battery's voltage, V , so each resistor in the circuit experiences a voltage drop $V/3$. Letting R stand for the resistance of one of the resistors, we find that the current through resistor B, which is the same as the currents through all the others, is given by $I_0 = V/3R$.

We now pass to the case where the switch is closed, as shown in the third panel. The battery's voltage is the same as before, and each resistor's resistance is the same, so we can still use the same symbols V and R for them. It is no longer true, however, that each resistor feels a voltage drop $V/3$. The equivalent resistance of the whole circuit is $R/2 + R/3 + R/2 = 4R/3$, so the total current drawn from the battery is $3V/4R$. In the middle group of resistors, this current is split three ways, so the new current through B is $(1/3)(3V/4R) = V/4R = 3I_0/4$.

Interpreting this result, we see that it comes from two effects that partially cancel. Closing the switch reduces the equivalent resistance of the circuit by giving charge another way to flow, and increases the amount of current drawn from the battery. Resistor B, however, only gets a 1/3 share of this greater current, not 1/2. The second effect turns out to be bigger than the first, and therefore the current through resistor B is lessened over all.

Discussion Question

A We have stated the loop rule in a symmetric form where a series of voltage drops adds up to zero. To do this, we had to define a standard way of connecting the voltmeter to the circuit so that the plus and minus signs would come out right. Suppose we wish to restate the junction rule in a similar symmetric way, so that instead of equating the current coming in to the current going out, it simply states that a certain sum of currents at a junction adds up to zero. What standard way of inserting the ammeter would we have to use to make this work?

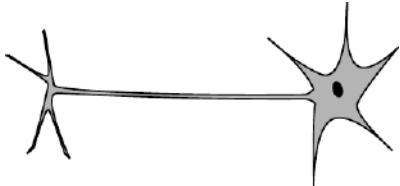
This chapter is summarized on page 1085. Notation and terminology are tabulated on pages 1070-1071.

Problems

The symbols \checkmark , \blacksquare , etc. are explained on page 573.

- 1** In a wire carrying a current of 1.0 pA , how long do you have to wait, on the average, for the next electron to pass a given point? Express your answer in units of microseconds.

\triangleright Solution, p. 1047 \blacksquare



Problem 2.

- 2** Referring back to our old friend the neuron from problem 1 on page 525, let's now consider what happens when the nerve is stimulated to transmit information. When the blob at the top (the cell body) is stimulated, it causes Na^+ ions to rush into the top of the tail (axon). This electrical pulse will then travel down the axon, like a flame burning down from the end of a fuse, with the Na^+ ions at each point first going out and then coming back in. If 10^{10} Na^+ ions cross the cell membrane in 0.5 ms , what amount of current is created?

\checkmark \blacksquare

- 3** If a typical light bulb draws about 900 mA from a 110 V household circuit, what is its resistance? (Don't worry about the fact that it's alternating current.) \checkmark \blacksquare

- 4** (a) Express the power dissipated by a resistor in terms of R and ΔV only, eliminating I . \checkmark

- (b) Electrical receptacles in your home are mostly 110 V , but circuits for electric stoves, air conditioners, and washers and driers are usually 220 V . The two types of circuits have differently shaped receptacles. Suppose you rewire the plug of a drier so that it can be plugged in to a 110 V receptacle. The resistor that forms the heating element of the drier would normally draw 200 W . How much power does it actually draw now? \checkmark \blacksquare

- 5** Lightning discharges a cloud during an electrical storm. Suppose that the current in the lightning bolt varies with time as $I = bt$, where b is a constant. Find the cloud's charge as a function of time. \checkmark \blacksquare

- 6** A resistor has a voltage difference ΔV across it, causing a current I to flow.

- (a) Find an equation for the power it dissipates as heat in terms of the variables I and R only, eliminating ΔV . \checkmark

- (b) If an electrical line coming to your house is to carry a given amount of current, interpret your equation from part a to explain whether the wire's resistance should be small, or large. \blacksquare

- 7** In AM (amplitude-modulated) radio, an audio signal $f(t)$ is multiplied by a sine wave $\sin \omega t$ in the megahertz frequency range. For simplicity, let's imagine that the transmitting antenna is a whip, and that charge goes back and forth between the top and bottom. Suppose that, during a certain time interval, the audio signal varies

linearly with time, giving a charge $q = (a + bt) \sin \omega t$ at the top of the whip and $-q$ at the bottom. Find the current as a function of time. ✓ ■

8 Problem 8 has been deleted. ■

9 Use the result of problem 38 on page 572 to find an equation for the voltage at a point in space at a distance r from a point charge Q . (Take your $V = 0$ distance to be anywhere you like.) ✓ ■

10 Hybrid and electric cars have been gradually gaining market share, but during the same period of time, manufacturers such as Porsche have also begun designing and selling cars with “mild hybrid” systems, in which power-hungry parts like water pumps are powered by a higher-voltage battery rather than running directly on shafts from the motor. Traditionally, car batteries have been 12 volts. Car companies have dithered over what voltage to use as the standard for mild hybrids, building systems based on 36 V, 42 V, and 48 V. For the purposes of this problem, we consider 36 V.

(a) Suppose the battery in a new car is used to run a device that requires the same amount of power as the corresponding device in the old car. Based on the sample figures above, how would the currents handled by the wires in one of the new cars compare with the currents in the old ones? ✓

(b) The real purpose of the greater voltage is to handle devices that need *more* power. Can you guess why they decided to change to higher-voltage batteries rather than increasing the power without increasing the voltage?



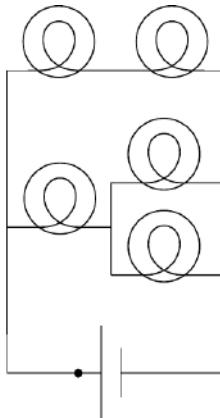
11 We have referred to resistors *dissipating* heat, i.e., we have assumed that $P = I\Delta V$ is always greater than zero. Could $I\Delta V$ come out to be negative for a resistor? If so, could one make a refrigerator by hooking up a resistor in such a way that it absorbed heat instead of dissipating it? ■

12 What resistance values can be created by combining a $1\text{ k}\Omega$ resistor and a $10\text{ k}\Omega$ resistor? ▷ Solution, p. 1047 ■

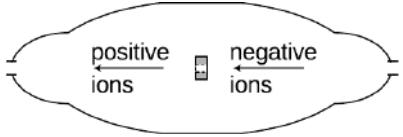
13 The figure shows a circuit containing five lightbulbs connected to a battery. Suppose you’re going to connect one probe of a voltmeter to the circuit at the point marked with a dot. How many unique, nonzero voltage differences could you measure by connecting the other probe to other wires in the circuit? ■

14 The lightbulbs in the figure are all identical. If you were inserting an ammeter at various places in the circuit, how many unique currents could you measure? If you know that the current measurement will give the same number in more than one place, only count that as one unique current. ■

15 (a) You take an LP record out of its sleeve, and it acquires a



Problems 13 and 14.



Problem 16.

static charge of 1 nC . You play it at the normal speed of $33\frac{1}{3} \text{ r.p.m.}$, and the charge moving in a circle creates an electric current. What is the current, in amperes? ✓

(b) Although the planetary model of the atom can be made to work with any value for the radius of the electrons' orbits, more advanced models that we will study later in this course predict definite radii. If the electron is imagined as circling around the proton at a speed of $2.2 \times 10^6 \text{ m/s}$, in an orbit with a radius of 0.05 nm , what electric current is created? ✓ ■

16 The figure shows a simplified diagram of a device called a tandem accelerator, used for accelerating beams of ions up to speeds on the order of 1-10% of the speed of light. (Since these velocities are not too big compared to c , you can use nonrelativistic physics throughout this problem.) The nuclei of these ions collide with the nuclei of atoms in a target, producing nuclear reactions for experiments studying the structure of nuclei. The outer shell of the accelerator is a conductor at zero voltage (i.e., the same voltage as the Earth). The electrode at the center, known as the “terminal,” is at a high positive voltage, perhaps millions of volts. Negative ions with a charge of -1 unit (i.e., atoms with one extra electron) are produced offstage on the right, typically by chemical reactions with cesium, which is a chemical element that has a strong tendency to give away electrons. Relatively weak electric and magnetic forces are used to transport these -1 ions into the accelerator, where they are attracted to the terminal. Although the center of the terminal has a hole in it to let the ions pass through, there is a very thin carbon foil there that they must physically penetrate. Passing through the foil strips off some number of electrons, changing the atom into a positive ion, with a charge of $+n$ times the fundamental charge. Now that the atom is positive, it is repelled by the terminal, and accelerates some more on its way out of the accelerator.

(a) Find the velocity, v , of the emerging beam of positive ions, in terms of n , their mass m , the terminal voltage V , and fundamental constants. Neglect the small change in mass caused by the loss of electrons in the stripper foil. ✓

(b) To fuse protons with protons, a minimum beam velocity of about 11% of the speed of light is required. What terminal voltage would be needed in this case? ✓

(c) In the setup described in part b, we need a target containing atoms whose nuclei are single protons, i.e., a target made of hydrogen. Since hydrogen is a gas, and we want a foil for our target, we have to use a hydrogen compound, such as a plastic. Discuss what effect this would have on the experiment. ■

17 Wire is sold in a series of standard diameters, called “gauges.”

The difference in diameter between one gauge and the next in the series is about 20%. How would the resistance of a given length of wire compare with the resistance of the same length of wire in the next gauge in the series? ✓ ■

18 In the figure, the battery is 9 V.

- (a) What are the voltage differences across each light bulb? ✓
- (b) What current flows through each of the three components of the circuit? ✓
- (c) If a new wire is added to connect points A and B, how will the appearances of the bulbs change? What will be the new voltages and currents? ■
- (d) Suppose no wire is connected from A to B, but the two bulbs are switched. How will the results compare with the results from the original setup as drawn? ■

19 A student in a biology lab is given the following instructions: "Connect the cerebral eraser (C.E.) and the neural depolarizer (N.D.) in parallel with the power supply (P.S.). (Under no circumstances should you ever allow the cerebral eraser to come within 20 cm of your head.) Connect a voltmeter to measure the voltage across the cerebral eraser, and also insert an ammeter in the circuit so that you can make sure you don't put more than 100 mA through the neural depolarizer." The diagrams show two lab groups' attempts to follow the instructions.

- (a) Translate diagram 1 into a standard-style schematic. What is correct and incorrect about this group's setup?
- (b) Do the same for diagram 2. ■

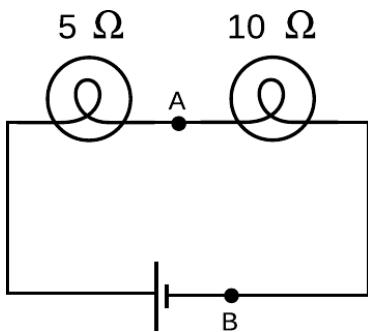
20 Referring back to problem 15, p. 527, about the sodium chloride crystal, suppose the lithium ion is going to jump from the gap it is occupying to one of the four closest neighboring gaps. Which one will it jump to, and if it starts from rest, how fast will it be going by the time it gets there? (It will keep on moving and accelerating after that, but that does not concern us.) You will need the result of problem 38.

▷ Hint, p. 1036 ✓ ■

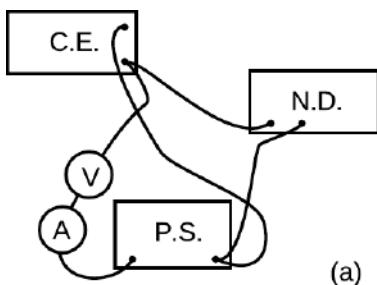
21 A $1.0\ \Omega$ toaster and a $2.0\ \Omega$ lamp are connected in parallel with the 110-V supply of your house. (Ignore the fact that the voltage is AC rather than DC.)

- (a) Draw a schematic of the circuit.
- (b) For each of the three components in the circuit, find the current passing through it and the voltage drop across it. ✓
- (c) Suppose they were instead hooked up in series. Draw a schematic and calculate the same things. ✓ ■

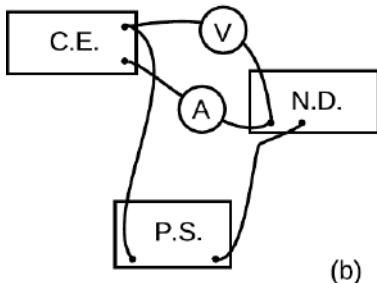
22 The heating element of an electric stove is connected in series with a switch that opens and closes many times per second. When you turn the knob up for more power, the fraction of the time that



Problem 18.

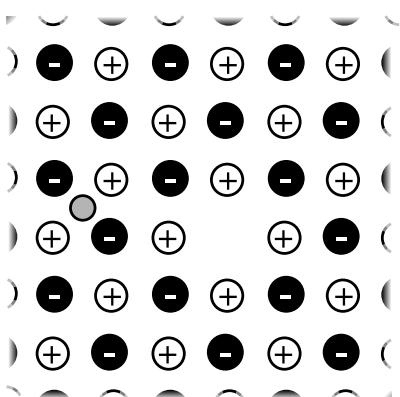


(a)



(b)

Problem 19.



Problem 20.

the switch is closed increases. Suppose someone suggests a simpler alternative for controlling the power by putting the heating element in series with a variable resistor controlled by the knob. (With the knob turned all the way clockwise, the variable resistor's resistance is nearly zero, and when it's all the way counterclockwise, its resistance is essentially infinite.) (a) Draw schematics. (b) Why would the simpler design be undesirable? ■

23 You have a circuit consisting of two unknown resistors in series, and a second circuit consisting of two unknown resistors in parallel.

- (a) What, if anything, would you learn about the resistors in the series circuit by finding that the currents through them were equal?
- (b) What if you found out the voltage differences across the resistors in the series circuit were equal?
- (c) What would you learn about the resistors in the parallel circuit from knowing that the currents were equal?
- (d) What if the voltages in the parallel circuit were equal? ■

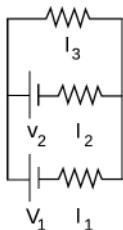
24 How many different resistance values can be created by combining three unequal resistors? (Don't count possibilities in which not all the resistors are used, i.e., ones in which there is zero current in one or more of them.) ■

25 Suppose six identical resistors, each with resistance R , are connected so that they form the edges of a tetrahedron (a pyramid with three sides in addition to the base, i.e., one less side than an Egyptian pyramid). What resistance value or values can be obtained by making connections onto any two points on this arrangement?

▷ Solution, p. 1047 ■

26 A person in a rural area who has no electricity runs an extremely long extension cord to a friend's house down the road so she can run an electric light. The cord is so long that its resistance, x , is not negligible. Show that the lamp's brightness is greatest if its resistance, y , is equal to x . Explain physically why the lamp is dim for values of y that are too small or too large. ■

27 All three resistors have the same resistance, R . Find the three unknown currents in terms of V_1 , V_2 , and R . ✓ ■



Problem 27.

28 You are given a battery, a flashlight bulb, and a single piece of wire. Draw at least two configurations of these items that would result in lighting up the bulb, and at least two that would not light it. (Don't draw schematics.) Note that the bulb has two electrical contacts: one is the threaded metal jacket, and the other is the tip (at the bottom in the figure).

If you're not sure what's going on, there are a couple of ways to check. The best is to try it in real life by either borrowing the materials from your instructor or scrounging the materials from around the house. (If you have a flashlight with this type of bulb, you can remove the bulb.) Another method is to use the simulation at phet.colorado.edu/en/simulation/circuit-construction-kit-dc.

[Problem by Arnold Arons.]

29 The figure shows a simplified diagram of an electron gun such as the one that creates the electron beam in a TV tube. Electrons that spontaneously emerge from the negative electrode (cathode) are then accelerated to the positive electrode, which has a hole in it. (Once they emerge through the hole, they will slow down. However, if the two electrodes are fairly close together, this slowing down is a small effect, because the attractive and repulsive forces experienced by the electron tend to cancel.)

(a) If the voltage difference between the electrodes is ΔV , what is the velocity of an electron as it emerges at B? Assume that its initial velocity, at A, is negligible, and that the velocity is nonrelativistic. (If you haven't read ch. 7 yet, don't worry about the remark about relativity.)

✓

(b) Evaluate your expression numerically for the case where $\Delta V=10$ kV, and compare to the speed of light. If you've read ch. 7 already, comment on whether the assumption of nonrelativistic motion was justified.

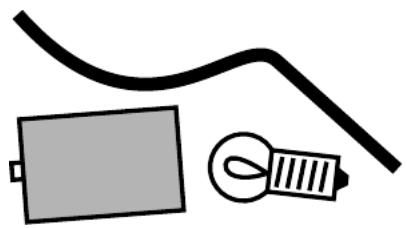
▷ Solution, p. 1048 ✓

30 (a) Many battery-operated devices take more than one battery. If you look closely in the battery compartment, you will see that the batteries are wired in series. Consider a flashlight circuit. What does the loop rule tell you about the effect of putting several batteries in series in this way?

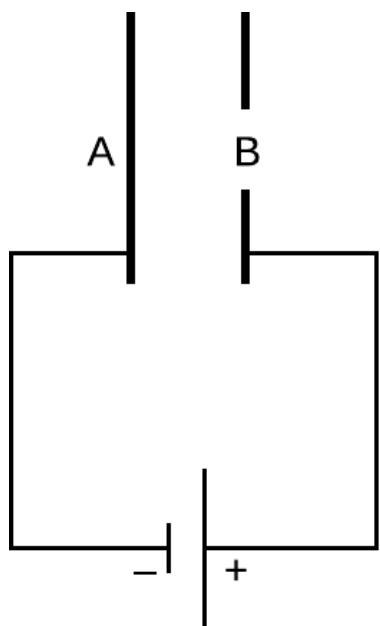
(b) The cells of an electric eel's nervous system are not that different from ours — each cell can develop a voltage difference across it of somewhere on the order of one volt. How, then, do you think an electric eel can create voltages of thousands of volts between different parts of its body?

■

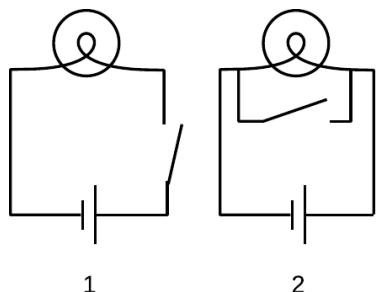
31 The figure shows two possible ways of wiring a flashlight with a switch. Both will serve to turn the bulb on and off, although the switch functions in the opposite sense. Why is method (1) preferable?



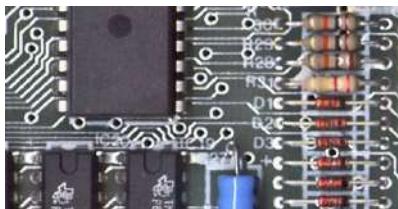
Problem 28.



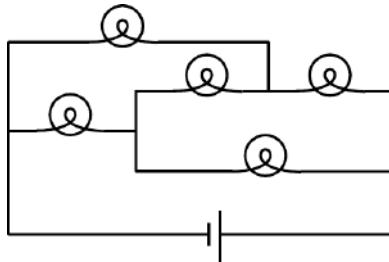
Problem 29.



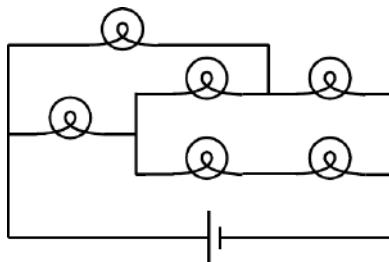
Problem 31.



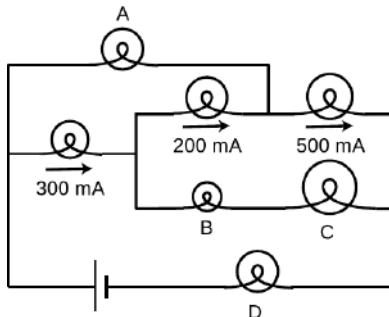
A printed circuit board, like the kind referred to in problem 32.



Problem 33.



Problem 34.



Problem 35.

32 You have to do different things with a circuit to measure current than to measure a voltage difference. Which would be more practical for a printed circuit board, in which the wires are actually strips of metal embedded inside the board?

► Solution, p. 1048 ■

33 The bulbs are all identical. Which one doesn't light up? ■

34 Each bulb has a resistance of one ohm. How much power is drawn from the one-volt battery? ✓ ■

35 The bulbs all have unequal resistances. Given the three currents shown in the figure, find the currents through bulbs A, B, C, and D. ✓ ■

36 A silk thread is uniformly charged by rubbing it with llama fur. The thread is then dangled vertically above a metal plate and released. As each part of the thread makes contact with the conducting plate, its charge is deposited onto the plate. Since the thread is accelerating due to gravity, the rate of charge deposition increases with time, and by time t the cumulative amount of charge is $q = ct^2$, where c is a constant. (a) Find the current flowing onto the plate. ✓ (b) Suppose that the charge is immediately carried away through a resistance R . Find the power dissipated as heat. ✓ ■

37 In example 9 on p. 546, suppose that the larger sphere has radius a , the smaller one b . (a) Use the result of problem 9 to show that the ratio of the charges on the two spheres is $q_a/q_b = a/b$. (b) Show that the density of charge (charge per unit area) is the other way around: the charge density on the smaller sphere is *greater* than that on the larger sphere in the ratio a/b . ■

38 (a) Recall that the gravitational energy of two gravitationally interacting spheres is given by $PE = -Gm_1m_2/r$, where r is the center-to-center distance. Sketch a graph of PE as a function of r , making sure that your graph behaves properly at small values of r , where you're dividing by a small number, and at large ones, where you're dividing by a large one. Check that your graph behaves properly when a rock is dropped from a larger r to a smaller one; the rock should *lose* potential energy as it gains kinetic energy. (b) Electrical forces are closely analogous to gravitational ones, since both depend on $1/r^2$. Since the forces are analogous, the potential energies should also behave analogously. Using this analogy, write down the expression for the electrical potential energy of two interacting charged particles. The main uncertainty here is the sign out in front. Like masses attract, but like charges repel. To figure out whether you have the right sign in your equation, sketch graphs in the case where both charges are positive, and also in the case where one is positive and one negative; make sure that in both cases, when the charges are released near one another, their motion causes them

to lose PE while gaining KE.

✓ ■

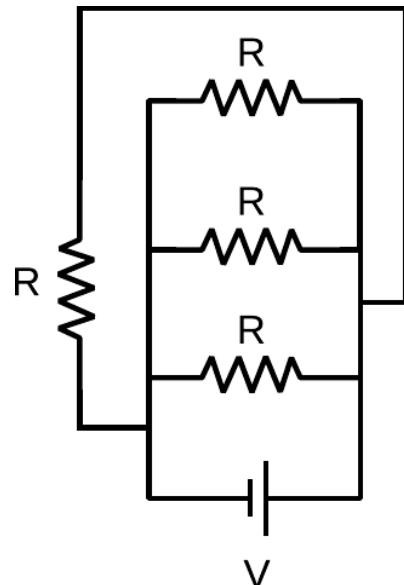
39 Find the current drawn from the battery.

✓ ■

40 It's fairly common in electrical circuits for additional, undesirable resistances to occur because of factors such as dirty, corroded, or loose connections. Suppose that a device with resistance R normally dissipates power P , but due to an additional series resistance r the *total* power is reduced to P' . We might, for example, detect this change because the battery powering our device ran down more slowly than normal.

- (a) Find the unknown resistance r .
- (b) Check that the units of your result make sense.
- (c) Check that your result makes sense in the special cases $P' = P$ and $P' = 0$.
- (d) Suppose we redefine P' as the useful power dissipated in R . For example, this would be the change we would notice because a flashlight was dimmer. Find r .

✓ ■



Problem 39.

Key to symbols:

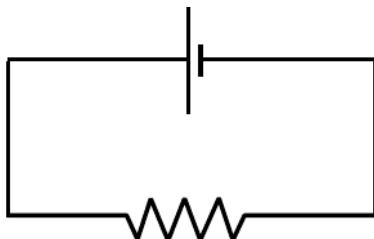
■ easy ■ typical ■ challenging ■ difficult ■ very difficult

✓ An answer check is available at www.lightandmatter.com.

Exercises

Exercise 9A: Voltage and Current

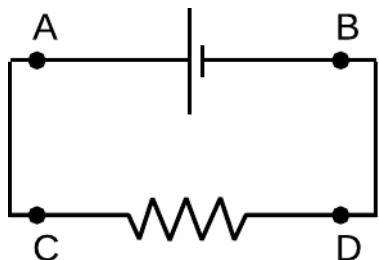
1. How many different currents could you measure in this circuit? Make a prediction, and then try it.



What do you notice? How does this make sense in terms of the roller coaster metaphor introduced in discussion question 9.1.3A on page 541?

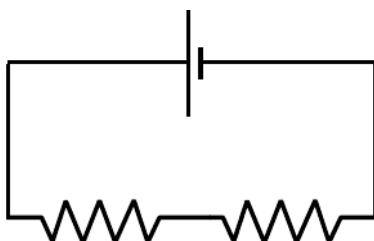
What is being *used up* in the resistor?

2. By connecting probes to these points, how many ways could you measure a voltage? How many of them would be different numbers? Make a prediction, and then do it.



What do you notice? Interpret this using the roller coaster metaphor, and color in parts of the circuit that represent constant voltages.

3. The resistors are unequal. How many *different* voltages and currents can you measure? Make a prediction, and then try it.



What do you notice? Interpret this using the roller coaster metaphor, and color in parts of the circuit that represent constant voltages.

Exercise 9B: The Loop and Junction Rules

Apparatus:

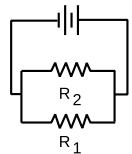
DC power supply

multimeter

resistors

1. The junction rule

Construct a circuit like this one, using the power supply as your voltage source. To make things more interesting, don't use equal resistors. Use nice big resistors (say $100\text{ k}\Omega$ to $1\text{ M}\Omega$) — this will ensure that you don't burn up the resistors, and that the multimeter's small internal resistance when used as an ammeter is negligible in comparison.



Insert your multimeter in the circuit to measure all three currents that you need in order to test the junction rule.

2. The loop rule

Now come up with a circuit to test the loop rule. Since the loop rule is always supposed to be true, it's hard to go wrong here! Make sure you have at least three resistors in a loop, and make sure you hook in the power supply in a way that creates non-zero voltage differences across all the resistors. Measure the voltage differences you need to measure to test the loop rule. Here it is best to use fairly small resistances, so that the multimeter's large internal resistance when used in parallel as a voltmeter will not significantly reduce the resistance of the circuit. Do not use resistances of less than about 100Ω , however, or you may blow a fuse or burn up a resistor.

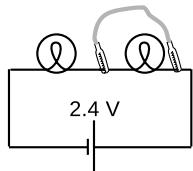
Exercise 9C: Reasoning About Circuits

The questions in this exercise can all be solved using some combination of the following approaches:

- There is constant voltage throughout any conductor.
- Ohm's law can be applied to any *part* of a circuit.
- Apply the loop rule.
- Apply the junction rule.

In each case, discuss the question, decide what you think is the right answer, and then try the experiment.

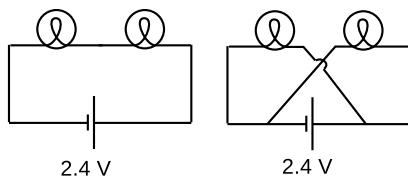
1. A wire is added in parallel with one bulb.



Which reasoning is correct?

- Each bulb still has 1.2 V across it, so both bulbs are still lit up.
- All parts of a wire are at the same voltage, and there is now a wire connection from one side of the right-hand bulb to the other. The right-hand bulb has no voltage difference across it, so it goes out.

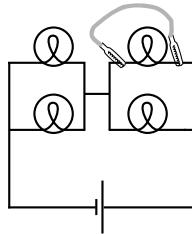
2. The series circuit is changed as shown.



Which reasoning is correct?

- Each bulb now has its sides connected to the two terminals of the battery, so each now has 2.4 V across it instead of 1.2 V. They get brighter.
- Just as in the original circuit, the current goes through one bulb, then the other. It's just that now the current goes in a figure-8 pattern. The bulbs glow the same as before.

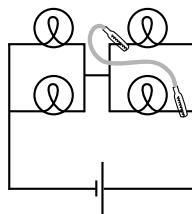
3. A wire is added as shown to the original circuit.



What is wrong with the following reasoning?

The top right bulb will go out, because its two sides are now connected with wire, so there will be no voltage difference across it. The other three bulbs will not be affected.

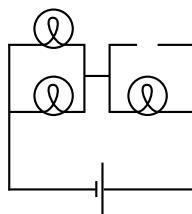
4. A wire is added as shown to the original circuit.



What is wrong with the following reasoning?

The current flows out of the right side of the battery. When it hits the first junction, some of it will go left and some will keep going up. The part that goes up lights the top right bulb. The part that turns left then follows the path of least resistance, going through the new wire instead of the bottom bulb. The top bulb stays lit, the bottom one goes out, and others stay the same.

5. What happens when one bulb is unscrewed, leaving an air gap?



Chapter 10

Fields

"Okay. Your duties are as follows: Get Breen. I don't care how you get him, but get him soon. That faker! He posed for twenty years as a scientist without ever being apprehended. Well, I'm going to do some apprehending that'll make all previous apprehending look like no apprehension at all. You with me?"

"Yes," said Battle, very much confused. "What's that thing you have?"

"Piggy-back heat-ray. You transpose the air in its path into an unstable isotope which tends to carry all energy as heat. Then you shoot your juice light, or whatever along the isotopic path and you burn whatever's on the receiving end. You want a few?"

"No," said Battle. "I have my gats. What else have you got for offense and defense?" Underbottam opened a cabinet and proudly waved an arm. "Everything," he said.

"Disintegrators, heat-rays, bombs of every type. And impenetrable shields of energy, massive and portable. What more do I need?"

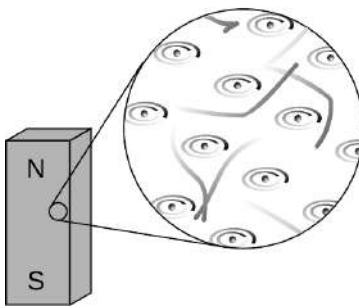
From THE REVERSIBLE REVOLUTIONS by Cecil Corwin, Cosmic Stories, March 1941. Art by Morey, Bok, Kyle, Hunt, Forte. Copyright expired.



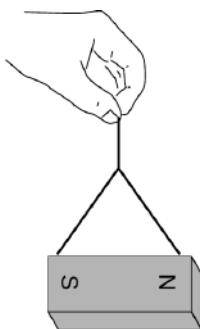
10.1 Fields of force

Cutting-edge science readily infiltrates popular culture, though sometimes in garbled form. The Newtonian imagination populated the universe mostly with that nice solid stuff called matter, which was made of little hard balls called atoms. In the early twentieth century, consumers of pulp fiction and popularized science began to hear of a new image of the universe, full of x-rays, N-rays, and Hertzian waves. What they were beginning to soak up through their skins was a drastic revision of Newton's concept of a universe made of chunks of matter which happened to interact via forces. In the newly emerging picture, the universe was *made* of force, or, to be more technically accurate, of ripples in universal fields of force. Unlike the average reader of Cosmic Stories in 1941, you now possess enough technical background to understand what a "force field" really is.

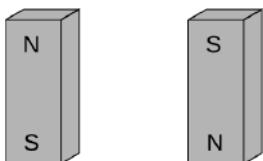
10.1.1 Why fields?



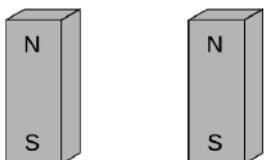
a / A bar magnet's atoms are (partially) aligned.



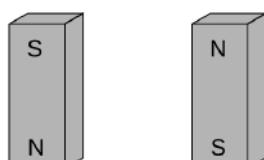
b / A bar magnet interacts with our magnetic planet.



c / Magnets aligned north-south.



d / The second magnet is reversed.



e / Both magnets are reversed.

Time delays in forces exerted at a distance

What convinced physicists that they needed this new concept of a field of force? Although we have been dealing mostly with electrical forces, let's start with a magnetic example. (In fact the main reason I've delayed a detailed discussion of magnetism for so long is that mathematical calculations of magnetic effects are handled much more easily with the concept of a field of force.) First a little background leading up to our example. A bar magnet, a, has an axis about which many of the electrons' orbits are oriented. The earth itself is also a magnet, although not a bar-shaped one. The interaction between the earth-magnet and the bar magnet, b, makes them want to line up their axes in opposing directions (in other words such that their electrons rotate in parallel planes, but with one set rotating clockwise and the other counterclockwise as seen looking along the axes). On a smaller scale, any two bar magnets placed near each other will try to align themselves head-to-tail, c.

Now we get to the relevant example. It is clear that two people separated by a paper-thin wall could use a pair of bar magnets to signal to each other. Each person would feel her own magnet trying to twist around in response to any rotation performed by the other person's magnet. The practical range of communication would be very short for this setup, but a sensitive electrical apparatus could pick up magnetic signals from much farther away. In fact, this is not so different from what a radio does: the electrons racing up and down the transmitting antenna create forces on the electrons in the distant receiving antenna. (Both magnetic and electric forces are involved in real radio signals, but we don't need to worry about that yet.)

A question now naturally arises as to whether there is any time delay in this kind of communication via magnetic (and electric) forces. Newton would have thought not, since he conceived of physics in terms of instantaneous action at a distance. We now know, however, that there is such a time delay. If you make a long-distance phone call that is routed through a communications satellite, you should easily be able to detect a delay of about half a second over the signal's round trip of 50,000 miles. Modern measurements have shown that electric, magnetic, and gravitational forces all travel at the speed of light, 3×10^8 m/s. (In fact, we will soon discuss how light itself is made of electricity and magnetism.)

If it takes some time for forces to be transmitted through space, then apparently there is some *thing* that travels *through* space. The fact that the phenomenon travels outward at the same speed in all directions strongly evokes wave metaphors such as ripples on a pond.

More evidence that fields of force are real: they carry energy.

The smoking-gun argument for this strange notion of traveling force ripples comes from the fact that they carry energy.

First suppose that the person holding the bar magnet on the right decides to reverse hers, resulting in configuration d. She had to do mechanical work to twist it, and if she releases the magnet, energy will be released as it flips back to c. She has apparently stored energy by going from c to d. So far everything is easily explained without the concept of a field of force.

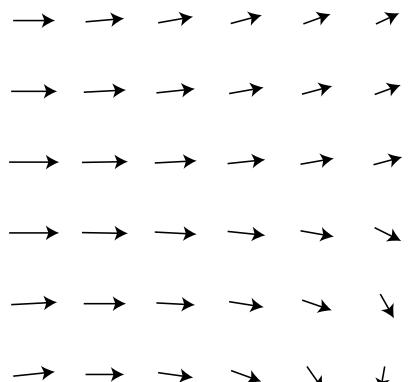
But now imagine that the two people start in position c and then simultaneously flip their magnets extremely quickly to position e, keeping them lined up with each other the whole time. Imagine, for the sake of argument, that they can do this so quickly that each magnet is reversed while the force signal from the other is still in transit. (For a more realistic example, we'd have to have two radio antennas, not two magnets, but the magnets are easier to visualize.) During the flipping, each magnet is still feeling the forces arising from the way the other magnet *used* to be oriented. Even though the two magnets stay aligned during the flip, the time delay causes each person to feel resistance as she twists her magnet around. How can this be? Both of them are apparently doing mechanical work, so they must be storing magnetic energy somehow. But in the traditional Newtonian conception of matter interacting via instantaneous forces at a distance, interaction energy arises from the relative positions of objects that are interacting via forces. If the magnets never changed their orientations relative to each other, how can any magnetic energy have been stored?

The only possible answer is that the energy must have gone into the magnetic force ripples crisscrossing the space between the magnets. Fields of force apparently carry energy across space, which is strong evidence that they are real things.

This is perhaps not as radical an idea to us as it was to our ancestors. We are used to the idea that a radio transmitting antenna consumes a great deal of power, and somehow spews it out into the universe. A person working around such an antenna needs to be careful not to get too close to it, since all that energy can easily cook flesh (a painful phenomenon known as an “RF burn”).

10.1.2 The gravitational field

Given that fields of force are real, how do we define, measure, and calculate them? A fruitful metaphor will be the wind patterns experienced by a sailing ship. Wherever the ship goes, it will feel a certain amount of force from the wind, and that force will be in a certain direction. The weather is ever-changing, of course, but for now let's just imagine steady wind patterns. Definitions in physics are operational, i.e., they describe how to measure the thing being



f / The wind patterns in a certain area of the ocean could be charted in a “sea of arrows” representation like this. Each arrow represents both the wind’s strength and its direction at a certain location.

defined. The ship's captain can measure the wind's "field of force" by going to the location of interest and determining both the direction of the wind and the strength with which it is blowing. Charting all these measurements on a map leads to a depiction of the field of wind force like the one shown in the figure. This is known as the "sea of arrows" method of visualizing a field.

Now let's see how these concepts are applied to the fundamental force fields of the universe. We'll start with the gravitational field, which is the easiest to understand. As with the wind patterns, we'll start by imagining gravity as a static field, even though the existence of the tides proves that there are continual changes in the gravity field in our region of space. When the gravitational field was introduced in chapter 2, I avoided discussing its direction explicitly, but defining it is easy enough: we simply go to the location of interest and measure the direction of the gravitational force on an object, such as a weight tied to the end of a string.

In chapter 2, I defined the gravitational field in terms of the energy required to raise a unit mass through a unit distance. However, I'm going to give a different definition now, using an approach that will be more easily adapted to electric and magnetic fields. This approach is based on force rather than energy. We couldn't carry out the energy-based definition without dividing by the mass of the object involved, and the same is true for the force-based definition. For example, gravitational forces are weaker on the moon than on the earth, but we cannot specify the strength of gravity simply by giving a certain number of newtons. The number of newtons of gravitational force depends not just on the strength of the local gravitational field but also on the mass of the object on which we're testing gravity, our "test mass." A boulder on the moon feels a stronger gravitational force than a pebble on the earth. We can get around this problem by defining the strength of the gravitational field as the force acting on an object, *divided by the object's mass*:

The gravitational field vector, \mathbf{g} , at any location in space is found by placing a test mass m_t at that point. The field vector is then given by $\mathbf{g} = \mathbf{F}/m_t$, where \mathbf{F} is the gravitational force on the test mass.

We now have three ways of representing a gravitational field. The magnitude of the gravitational field near the surface of the earth, for instance, could be written as 9.8 N/kg, 9.8 J/kg · m, or 9.8 m/s². If we already had two names for it, why invent a third? The main reason is that it prepares us with the right approach for defining other fields.

The most subtle point about all this is that the gravitational field tells us about what forces *would* be exerted on a test mass by the earth, sun, moon, and the rest of the universe, *if* we inserted a test mass at the point in question. The field still exists at all the

places where we didn't measure it.

Gravitational field of the earth

example 1

- ▷ What is the magnitude of the earth's gravitational field, in terms of its mass, M , and the distance r from its center?
- ▷ Substituting $|\mathbf{F}| = GMm_t/r^2$ into the definition of the gravitational field, we find $|\mathbf{g}| = GM/r^2$. This expression could be used for the field of any spherically symmetric mass distribution, since the equation we assumed for the gravitational force would apply in any such case.

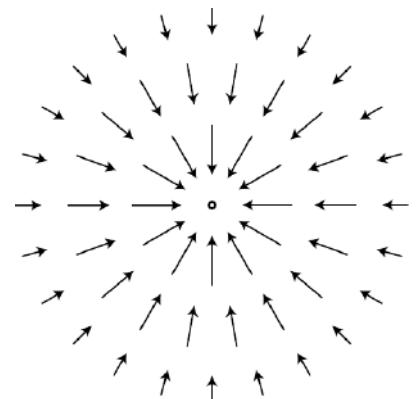
Sources and sinks

If we make a sea-of-arrows picture of the gravitational fields surrounding the earth, \mathbf{g} , the result is evocative of water going down a drain. For this reason, anything that creates an inward-pointing field around itself is called a sink. The earth is a gravitational sink. The term "source" can refer specifically to things that make outward fields, or it can be used as a more general term for both "outies" and "innies." However confusing the terminology, we know that gravitational fields are only attractive, so we will never find a region of space with an outward-pointing field pattern.

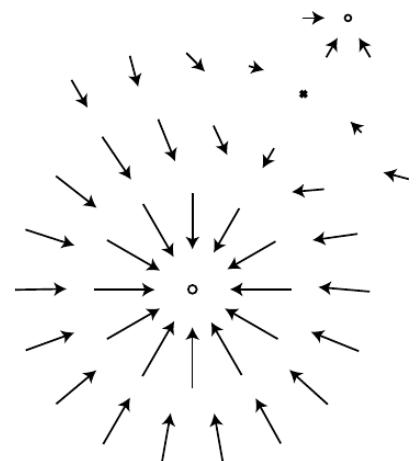
Knowledge of the field is interchangeable with knowledge of its sources (at least in the case of a static, unchanging field). If aliens saw the earth's gravitational field pattern they could immediately infer the existence of the planet, and conversely if they knew the mass of the earth they could predict its influence on the surrounding gravitational field.

Superposition of fields

A very important fact about all fields of force is that when there is more than one source (or sink), the fields add according to the rules of vector addition. The gravitational field certainly will have this property, since it is defined in terms of the force on a test mass, and forces add like vectors. Superposition is an important characteristic of waves, so the superposition property of fields is consistent with the idea that disturbances can propagate outward as waves in a field.



g / The gravitational field surrounding a clump of mass such as the earth.



h / The gravitational fields of the earth and moon superpose. Note how the fields cancel at one point, and how there is no boundary between the interpenetrating fields surrounding the two bodies.

Reduction in gravity on Io due to Jupiter's gravity example 2

▷ The average gravitational field on Jupiter's moon Io is 1.81 N/kg. By how much is this reduced when Jupiter is directly overhead? Io's orbit has a radius of 4.22×10^8 m, and Jupiter's mass is 1.899×10^{27} kg.

▷ By the shell theorem, we can treat the Jupiter as if its mass was all concentrated at its center, and likewise for Io. If we visit Io and land at the point where Jupiter is overhead, we are on the same line as these two centers, so the whole problem can be treated one-dimensionally, and vector addition is just like scalar addition. Let's use positive numbers for downward fields (toward the center of Io) and negative for upward ones. Plugging the appropriate data into the expression derived in example 1, we find that the Jupiter's contribution to the field is -0.71 N/kg. Superposition says that we can find the actual gravitational field by adding up the fields created by Io and Jupiter: $1.81 - 0.71$ N/kg = 1.1 N/kg. You might think that this reduction would create some spectacular effects, and make Io an exciting tourist destination. Actually you would not detect any difference if you flew from one side of Io to the other. This is because your body and Io both experience Jupiter's gravity, so you follow the same orbital curve through the space around Jupiter.



i / The part of the LIGO gravity wave detector at Hanford Nuclear Reservation, near Richland, Washington. The other half of the detector is in Louisiana.

Gravitational waves

A source that sits still will create a static field pattern, like a steel ball sitting peacefully on a sheet of rubber. A moving source will create a spreading wave pattern in the field, like a bug thrash-

ing on the surface of a pond. Although we have started with the gravitational field as the simplest example of a static field, stars and planets do more stately gliding than thrashing, so gravitational waves are not easy to detect. Newton's theory of gravity does not describe gravitational waves, but they are predicted by Einstein's general theory of relativity.

A Caltech-MIT collaboration has built a pair of gravitational wave detectors called LIGO to search for direct evidence of gravitational waves. Since they are essentially the most sensitive vibration detectors ever made, they are located in quiet rural areas, and signals are compared between them to make sure that they were not due to passing trucks. The signature of a gravitational wave is if the same wiggle is seen in both detectors within a short time. The detectors are able to sense a vibration that causes a change of 10^{-18} m in the distance between the mirrors at the ends of the 4-km vacuum tunnels. This is a thousand times less than the size of an atomic nucleus! In 2016, the collaboration announced the first detection of a gravitational wave, which is believed to have originated from the collision of two black holes. Propagation of gravitational waves at c was verified through multiple methods both by study of the 2016 event and through an event in 2017, interpreted as a collision of two neutron stars, in which both gravitational waves and electromagnetic waves were detected simultaneously.

10.1.3 The electric field

Definition

The definition of the electric field is directly analogous to, and has the same motivation as, the definition of the gravitational field:

The electric field vector, \mathbf{E} , at any location in space is found by placing a test charge q_t at that point. The electric field vector is then given by $\mathbf{E} = \mathbf{F}/q_t$, where \mathbf{F} is the electric force on the test charge.

Charges are what create electric fields. Unlike gravity, which is always attractive, electricity displays both attraction and repulsion. A positive charge is a source of electric fields, and a negative one is a sink.

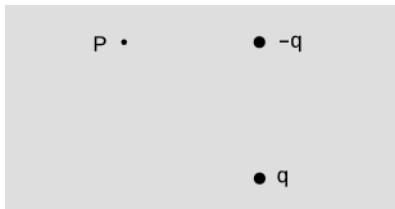
The most difficult point about the definition of the electric field is that the force on a negative charge is in the opposite direction compared to the field. This follows from the definition, since dividing a vector by a negative number reverses its direction. It's as though we had some objects that fell upward instead of down.

self-check A

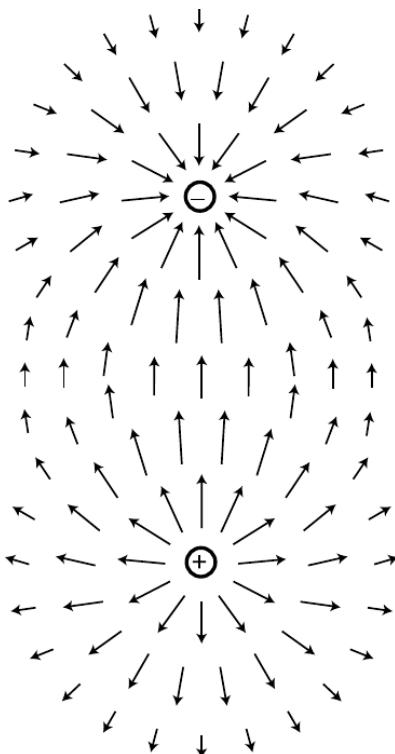
Find an equation for the magnitude of the field of a single point charge Q . ▷ Answer, p. 1063

Superposition of electric fields

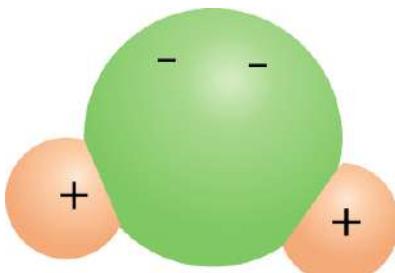
example 3



j / Example 3.



k / A dipole field. Electric fields diverge from a positive charge and converge on a negative charge.



l / A water molecule is a dipole.

▷ Charges q and $-q$ are at a distance b from each other, as shown in the figure. What is the electric field at the point P, which lies at a third corner of the square?

▷ The field at P is the vector sum of the fields that would have been created by the two charges independently. Let positive x be to the right and let positive y be up.

Negative charges have fields that point at them, so the charge $-q$ makes a field that points to the right, i.e., has a positive x component. Using the answer to the self-check, we have

$$E_{-q,x} = \frac{kq}{b^2}$$

$$E_{-q,y} = 0.$$

Note that if we had blindly ignored the absolute value signs and plugged in $-q$ to the equation, we would have incorrectly concluded that the field went to the left.

By the Pythagorean theorem, the positive charge is at a distance $\sqrt{2}b$ from P, so the magnitude of its contribution to the field is $E = kq/2b^2$. Positive charges have fields that point away from them, so the field vector is at an angle of 135° counterclockwise from the x axis.

$$E_{q,x} = \frac{kq}{2b^2} \cos 135^\circ$$

$$= -\frac{kq}{2^{3/2}b^2}$$

$$E_{q,y} = \frac{kq}{2b^2} \sin 135^\circ$$

$$= \frac{kq}{2^{3/2}b^2}$$

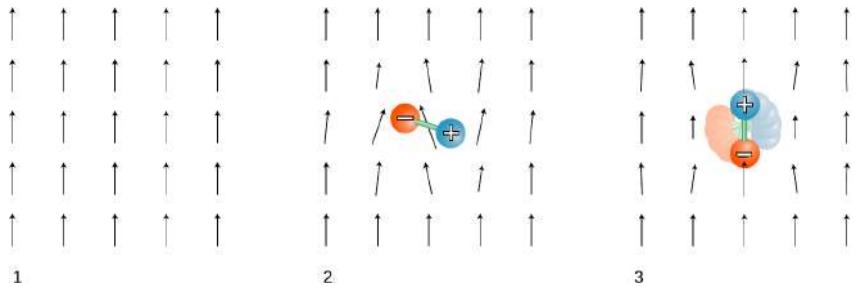
The total field is

$$E_x = \left(1 - 2^{-3/2}\right) \frac{kq}{b^2}$$

$$E_y = \frac{kq}{2^{3/2}b^2}$$

Dipoles

The simplest set of sources that can occur with electricity but not with gravity is the *dipole*, consisting of a positive charge and a negative charge with equal magnitudes. More generally, an electric dipole can be any object with an imbalance of positive charge on one side and negative on the other. A water molecule, l, is a dipole because the electrons tend to shift away from the hydrogen atoms and onto the oxygen atom.



- m / 1. A uniform electric field created by some charges “off-stage.”
 2. A dipole is placed in the field. 3. The dipole aligns with the field.

Your microwave oven acts on water molecules with electric fields. Let us imagine what happens if we start with a uniform electric field, m/1, made by some external charges, and then insert a dipole, m/2, consisting of two charges connected by a rigid rod. The dipole disturbs the field pattern, but more important for our present purposes is that it experiences a torque. In this example, the positive charge feels an upward force, but the negative charge is pulled down. The result is that the dipole wants to align itself with the field, m/3. The microwave oven heats food with electrical (and magnetic) waves. The alternation of the torque causes the molecules to wiggle and increase the amount of random motion. The slightly vague definition of a dipole given above can be improved by saying that a dipole is any object that experiences a torque in an electric field.

What determines the torque on a dipole placed in an externally created field? Torque depends on the force, the distance from the axis at which the force is applied, and the angle between the force and the line from the axis to the point of application. Let a dipole consisting of charges $+q$ and $-q$ separated by a distance ℓ be placed in an external field of magnitude $|\mathbf{E}|$, at an angle θ with respect to the field. The total torque on the dipole is

$$\begin{aligned}\tau &= \frac{\ell}{2}q|\mathbf{E}|\sin\theta + \frac{\ell}{2}q|\mathbf{E}|\sin\theta \\ &= \ell q|\mathbf{E}|\sin\theta.\end{aligned}$$

(Note that even though the two forces are in opposite directions, the torques do not cancel, because they are both trying to twist the dipole in the same direction.) The quantity is called the dipole moment, notated D . (More complex dipoles can also be assigned a dipole moment — they are defined as having the same dipole moment as the two-charge dipole that would experience the same torque.)

Employing a little more mathematical elegance, we can define a dipole moment *vector*,

$$\mathbf{D} = \sum q_i \mathbf{r}_i,$$

where \mathbf{r}_i is the position vector of the charge labeled by the index i . We can then write the torque in terms of a vector cross product (page 287),

$$\boldsymbol{\tau} = \mathbf{D} \times \mathbf{E}.$$

No matter how we notate it, the definition of the dipole moment requires that we choose a point from which we measure all the position vectors of the charges. However, in the commonly encountered special case where the total charge of the object is zero, the dipole moment is the same regardless of this choice.

Dipole moment of a molecule of NaCl gas

example 4

▷ In a molecule of NaCl gas, the center-to-center distance between the two atoms is about 0.24 nm. Assuming that the chlorine completely steals one of the sodium's electrons, compute the magnitude of this molecule's dipole moment.

▷ The total charge is zero, so it doesn't matter where we choose the origin of our coordinate system. For convenience, let's choose it to be at one of the atoms, so that the charge on that atom doesn't contribute to the dipole moment. The magnitude of the dipole moment is then

$$\begin{aligned} D &= (2.4 \times 10^{-10} \text{ m})(e) \\ &= (2.4 \times 10^{-10} \text{ m})(1.6 \times 10^{-19} \text{ C}) \\ &\approx 4 \times 10^{-29} \text{ C} \cdot \text{m} \end{aligned}$$

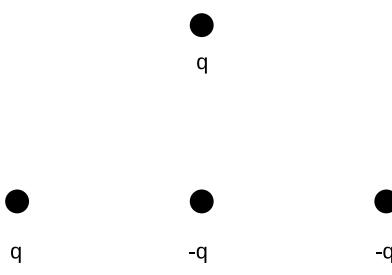
The experimentally measured value is $3.0 \times 10^{-29} \text{ C} \cdot \text{m}$, which shows that the electron is not completely "stolen."

Dipole moments as vectors

example 5

▷ The horizontal and vertical spacing between the charges in the figure is b . Find the dipole moment.

▷ Let the origin of the coordinate system be at the leftmost charge.



$$\begin{aligned} \mathbf{D} &= \sum q_i \mathbf{r}_i \\ &= (q)(0) + (-q)(b\hat{x}) + (q)(b\hat{x} + b\hat{y}) + (-q)(2b\hat{x}) \\ &= -2bq\hat{x} + bq\hat{y} \end{aligned}$$

n / Example 5.

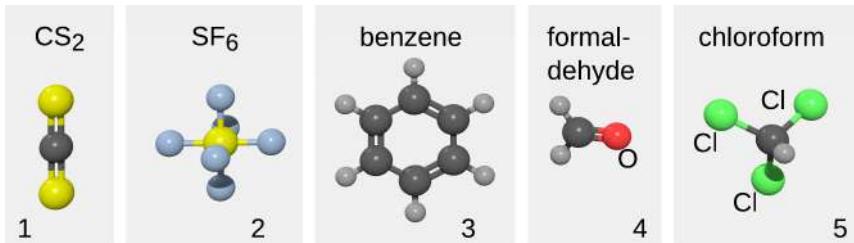
The dipole moment points up and to the left, which makes sense because the positive charges are predominantly above and to the left of the negative ones.

Molecules with zero and nonzero dipole moments

example 6

It can be useful to know whether or not a molecule is polar, i.e., has a nonzero dipole moment. A polar molecule such as water is readily heated in a microwave oven, while a nonpolar one is not. Polar molecules are attracted to one another, so polar substances dissolve in other polar substances, but not in nonpolar substances, i.e., "like dissolves like."

In a symmetric molecule such as carbon disulfide, figure o/1, the dipole moment vanishes by symmetry. For if we rotate the molecule by 180 degrees about any one of the three coordinate axes defined in the caption of the figure, the molecule is unchanged, which means that its dipole moment is unchanged. A zero vector is the only vector that can stay the same under all these rotations.



o / Example 6. The positive x axis is to the right, y is up, and z is out of the page. Dark gray atoms are carbon, and the small light gray ones are hydrogen. Some other elements are labeled when their identity would otherwise not be clear.

If we wish, we can think of this vanishing of the dipole moment as arising from a cancellation of two dipole vectors, one for each of the molecular bonds. Each chemical element has a certain affinity for electrons, which is represented by the column in which it lies on the periodic table, and is ultimately determined by quantum physics (section 13.4.7, p. 940). Sulfur “likes” electrons just slightly more than carbon, so that the electron cloud shifts outward somewhat from the center. This gives rise to two dipole vectors, each of which points inward, toward the positive charge. (Chemists use an opposite convention for the direction.) These two vectors add up to zero.

Similar symmetry arguments show that sulfur hexafluoride, o/2, and benzene o/3, have vanishing dipole moments.

The formaldehyde molecule, o/4, does not have enough symmetry to guarantee that its dipole moment must vanish, but it does have enough to dictate that the dipole vector must lie along the x axis. To determine the sign of D_x , we must use the fact that oxygen has a much higher electron affinity than carbon, so this creates a strong contribution to the dipole moment in the negative x direction. Carbon in turn has a higher electron affinity than hydrogen, so there is also a dipole moment associated with each CH bond, pointing toward the hydrogen and therefore contributing a further negative D_x . The y components cancel. A similar analysis for chloroform, o/5, shows that the dipole moment points in the positive z direction.

From these arguments we can tell, for example, that carbon disulfide will be soluble in benzene, but chloroform will not.

Alternative definition of the electric field

The behavior of a dipole in an externally created field leads us to an alternative definition of the electric field:

The electric field vector, E , at any location in space is defined by observing the torque exerted on a test dipole D_t placed there. The direction of the field is the direction in which the field tends to align a dipole (from $-$ to $+$), and the field's magnitude is $|E| = \tau/D_t \sin \theta$. In other words, the field vector is the vector that satisfies the equation $\boldsymbol{\tau} = \mathbf{D}_t \times \mathbf{E}$ for any test dipole \mathbf{D}_t placed at that point in space.

The main reason for introducing a second definition for the same concept is that the magnetic field is most easily defined using a similar approach.

Energy of a dipole in a field

A dipole in an external field has a stable equilibrium orientation in which \mathbf{D} is parallel to \mathbf{E} . Since the vector cross product vanishes for parallel vectors, the field's torque on the dipole is zero. (There is also an unstable equilibrium with \mathbf{D} and \mathbf{E} antiparallel.) This fact is probably more familiar from the magnetic context, where a magnetic dipole such as a compass needle tends to align itself with an external magnetic field. We will encounter magnetic fields and dipoles later, and they have many properties analogous to those of their electric counterparts. Depending on its orientation, the dipole will have some interaction energy U when it interacts with the field. If the physical size of the dipole is small, the electric field can be approximated as having a single value throughout. Since energy is a scalar, and the dot product is the only way (up to a multiplicative constant) to multiply two vectors to get a scalar (sec. 3.4.5, p. 216), we must have $U \propto \mathbf{D} \cdot \mathbf{E}$. Because the energy is minimized when \mathbf{D} and \mathbf{E} are parallel, the constant of proportionality must be negative, and one can easily show by considering a concrete example that this constant equals -1 . Therefore we have

$$U = -\mathbf{D} \cdot \mathbf{E}. \quad [\text{energy of a pointlike dipole in a field}]$$

Force on a dipole in a nonuniform field

example 7

If a dipole with zero total charge is placed in a uniform field, it may experience a torque, but it will not experience any *force*. The situation changes if the field is nonuniform. The force can be nonzero, and, perhaps more surprisingly, the force depends only on the dipole moment, not on the details of the arrangement of the charges inside the dipole, provided that the dipole is small enough. This can be shown either by a brute-force calculation (problem 59, 668) or by the following slightly slicker technique. The force in the x direction is

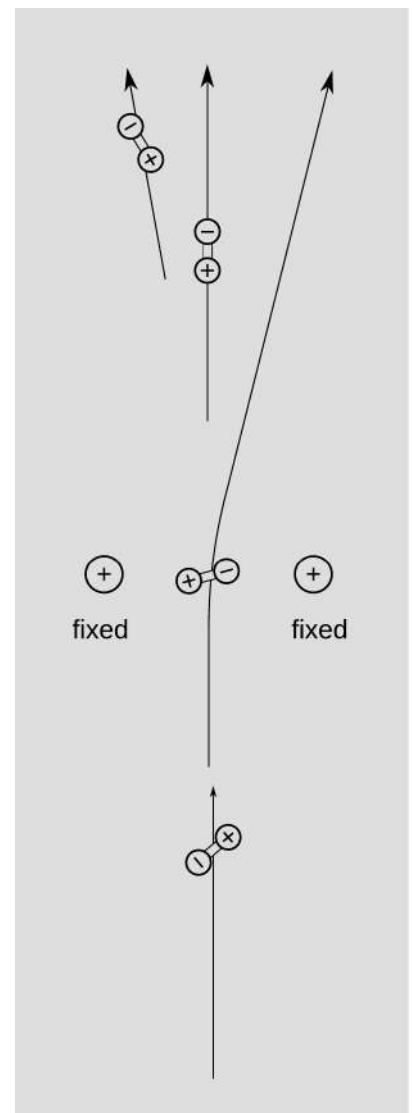
$$F_x = -\frac{\partial U}{\partial x},$$

where the symbol ∂ indicates a derivative with respect to the x coordinate only, holding y and z constant. (This is referred to as a partial derivative.) We then have

$$F_x = -\frac{\partial}{\partial x}(\mathbf{D} \cdot \mathbf{E}) = \mathbf{D} \cdot \frac{\partial \mathbf{E}}{\partial x},$$

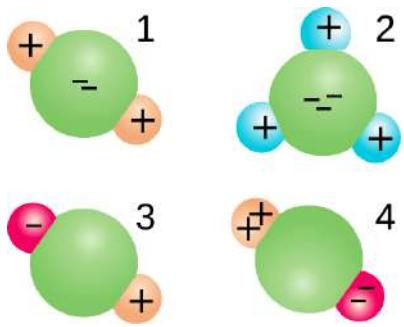
which depends on the dipole's properties only through \mathbf{D} . Similar expressions apply for the y and z components.

This principle can be used as a way of measuring the unknown dipole moments of a beam of particles, as in figure p. At a more pedestrian level, this is one way of explaining the fact that we can use a charged object to pick up uncharged scraps of paper (p. 478).



p / An electric dipole spectrometer. A beam of randomly oriented dipoles is shot through a "croquet hoop" consisting of two fixed positive charges. Although the field along the central axis of symmetry equals zero, the field is nonuniform, and therefore the dipoles feel a nonvanishing force, and are sorted out according to their orientations. A magnetic version of this device, described on p. 745, was used in the historic Stern-Gerlach experiment, sec. 14.1, p. 959.

Discussion Questions



q / Discussion

question

H.

A In the definition of the electric field, does the test charge need to be 1 coulomb? Does it need to be positive?

B Does a charged particle such as an electron or proton feel a force from its own electric field?

C Is there an electric field surrounding a wall socket that has nothing plugged into it, or a battery that is just sitting on a table?

D In a flashlight powered by a battery, which way do the electric fields point? What would the fields be like inside the wires? Inside the filament of the bulb?

E Criticize the following statement: "An electric field can be represented by a sea of arrows showing how current is flowing."

F The field of a point charge, $|\mathbf{E}| = kQ/r^2$, was derived in a self-check. How would the field pattern of a uniformly charged sphere compare with the field of a point charge?

G The interior of a perfect electrical conductor in equilibrium must have zero electric field, since otherwise the free charges within it would be drifting in response to the field, and it would not be in equilibrium. What about the field right at the surface of a perfect conductor? Consider the possibility of a field perpendicular to the surface or parallel to it.

H Compare the dipole moments of the molecules and molecular ions shown in the figure.

I Small pieces of paper that have not been electrically prepared in any way can be picked up with a charged object such as a charged piece of tape. In our new terminology, we could describe the tape's charge as inducing a dipole moment in the paper. Can a similar technique be used to induce not just a dipole moment but a charge?

10.2 Potential related to field

10.2.1 One dimension

Electrical potential (voltage) is electrical energy per unit charge, and electric field is force per unit charge. For a particle moving in one dimension, along the x axis, we can therefore relate potential and field if we start from the relationship between interaction energy and force,

$$dU = -F_x dx,$$

and divide by charge,

$$\frac{dU}{q} = -\frac{F_x}{q} dx,$$

giving

$$dV = -E_x dx,$$

or

$$\frac{dV}{dx} = -E_x.$$

The interpretation is that a strong electric field occurs in a region of space where the potential is rapidly changing. By analogy, a steep hillside is a place on the map where the altitude is rapidly changing.

Field generated by an electric eel

example 8

- ▷ Suppose an electric eel is 1 m long, and generates a voltage difference of 1000 volts between its head and tail. What is the electric field in the water around it?
- ▷ We are only calculating the amount of field, not its direction, so we ignore positive and negative signs. Subject to the possibly inaccurate assumption of a constant field parallel to the eel's body, we have

$$\begin{aligned} |\mathbf{E}| &= \frac{dV}{dx} \\ &\approx \frac{\Delta V}{\Delta x} \quad [\text{assumption of constant field}] \\ &= 1000 \text{ V/m}. \end{aligned}$$

Relating the units of electric field and potential

example 9

From our original definition of the electric field, we expect it to have units of newtons per coulomb, N/C. The example above, however, came out in volts per meter, V/m. Are these inconsistent? Let's reassure ourselves that this all works. In this kind of situation, the best strategy is usually to simplify the more complex units so that they involve only mks units and coulombs. Since potential is defined as electrical energy per unit charge, it has units of J/C:

$$\begin{aligned} \frac{V}{m} &= \frac{J/C}{m} \\ &= \frac{J}{C \cdot m}. \end{aligned}$$

To connect joules to newtons, we recall that work equals force times distance, so $J = N \cdot m$, so

$$\begin{aligned} \frac{V}{m} &= \frac{N \cdot m}{C \cdot m} \\ &= \frac{N}{C} \end{aligned}$$

As with other such difficulties with electrical units, one quickly begins to recognize frequently occurring combinations.

Potential associated with a point charge

example 10

- ▷ What is the potential associated with a point charge?
- ▷ As derived previously in self-check A on page 585, the field is

$$|\mathbf{E}| = \frac{kQ}{r^2}$$

The difference in potential between two points on the same radius line is

$$\begin{aligned}\Delta V &= - \int dV \\ &= - \int E_x dx\end{aligned}$$

In the general discussion above, x was just a generic name for distance traveled along the line from one point to the other, so in this case x really means r .

$$\begin{aligned}\Delta V &= - \int_{r_1}^{r_2} E_r dr \\ &= - \int_{r_1}^{r_2} \frac{kQ}{r^2} dr \\ &= \left[\frac{kQ}{r} \right]_{r_1}^{r_2} \\ &= \frac{kQ}{r_2} - \frac{kQ}{r_1}.\end{aligned}$$

The standard convention is to use $r_1 = \infty$ as a reference point, so that the potential at any distance r from the charge is

$$V = \frac{kQ}{r}.$$

The interpretation is that if you bring a positive test charge closer to a positive charge, its electrical energy is increased; if it was released, it would spring away, releasing this as kinetic energy.

self-check B

Show that you can recover the expression for the field of a point charge by evaluating the derivative $E_x = -dV/dx$. ▷ Answer, p. 1063

Weighing an electron

example 11

J.J. Thomson (p. 489) is considered to have discovered the electron because he measured its charge-to-mass ratio q/m and found it to be much larger than that of an ionized atom, interpreting this as evidence that he was seeing a subatomic particle with a mass much smaller than an atom's. But not only is the electron's q/m relatively large compared to that of an atom, it is simply a huge number ($\sim -10^{11}$ C/kg) when expressed in SI units. SI units are designed for human scales of experience, so this suggests that in

everyday life we should expect it to be very difficult to detect any effect from the weight or inertia of an electron.

As an example, suppose that a metal rod of length L is oriented upright. The conduction electrons are free to move, so they would tend to drop to the bottom of the rod. Electrical forces will however resist this segregation of positive and negative charges. To estimate how hard it would be to observe such an effect, let us imagine connecting the probes of a voltmeter to the ends of the rod. In equilibrium, the electrical and gravitational fields must have effects on an electron that cancel out. Setting the magnitudes of these forces equal to each other, we have $eE = mg$, and since $E = \Delta V/L$, we predict a voltage difference $\Delta V = (m/q)gL$. For a one-meter rod, the predicted effect is $\sim 10^{-10}$ V.

This is quite small, but not impossible to measure, and the theoretical prediction was confirmed for a similar experiment by Tolman and Stewart in a 1916 experiment at Berkeley. This was the first direct evidence that the charge carriers inside a metal wire are in fact electrons. Similarly, we do expect mechanical side-effects in any electrical circuit, e.g., a slight twitching of a flashlight when we turn it on or off, but these will be much too small to notice except with exceptionally delicate and sensitive tools. It is surprising that we can get information about the microscopic structure of a metal merely by measuring its properties in this way. Another, similar example along these lines is described in sec. 11.2.4, p. 697.

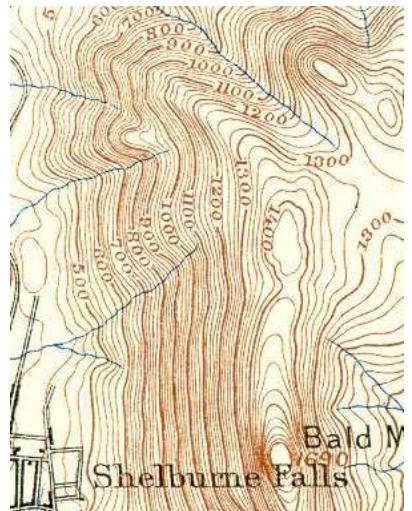
10.2.2 Two or three dimensions

The topographical map in figure a suggests a good way to visualize the relationship between field and potential in two dimensions. Each contour on the map is a line of constant height; some of these are labeled with their elevations in units of feet. Height is related to gravitational energy, so in a gravitational analogy, we can think of height as representing potential. Where the contour lines are far apart, as in the town, the slope is gentle. Lines close together indicate a steep slope.

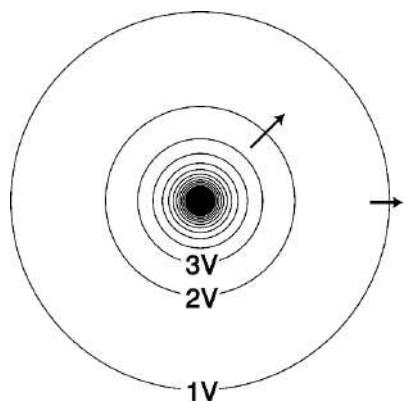
If we walk along a straight line, say straight east from the town, then height (potential) is a function of the east-west coordinate x . Using the usual mathematical definition of the slope, and writing V for the height in order to remind us of the electrical analogy, the slope along such a line is dV/dx (the rise over the run).

What if everything isn't confined to a straight line? Water flows downhill. Notice how the streams on the map cut perpendicularly through the lines of constant height.

It is possible to map potentials in the same way, as shown in figure b. The electric field is strongest where the constant-potential curves are closest together, and the electric field vectors always point



a / A topographical map of Shelburne Falls, Mass. (USGS)



b / The constant-potential curves surrounding a point charge. Near the charge, the curves are so closely spaced that they blend together on this drawing due to the finite width with which they were drawn. Some electric fields are shown as arrows.

perpendicular to the constant-potential curves.

The one-dimensional relationship $E = -dV/dx$ generalizes to three dimensions as follows:

$$E_x = -\frac{dV}{dx}$$
$$E_y = -\frac{dV}{dy}$$
$$E_z = -\frac{dV}{dz}$$

This can be notated as a gradient (page 219),

$$\mathbf{E} = -\nabla V,$$

and if we know the field and want to find the potential, we can use a line integral,

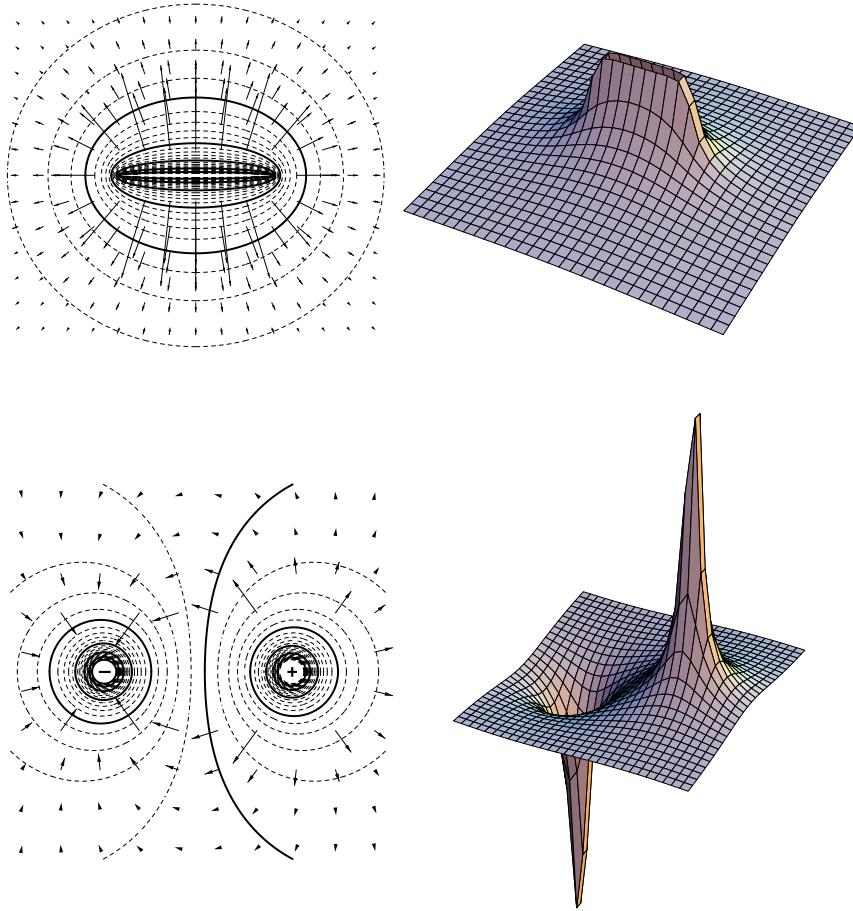
$$\Delta V = - \int_C \mathbf{E} \cdot d\mathbf{r},$$

where the quantity inside the integral is a vector dot product.

self-check C

Imagine that figure a represents potential rather than height. (a) Consider the stream the starts near the center of the map. Determine the positive and negative signs of dV/dx and dV/dy , and relate these to the direction of the force that is pushing the current forward against the resistance of friction. (b) If you wanted to find a lot of electric charge on this map, where would you look? ▷ Answer, p. 1063

Figure c shows some examples of ways to visualize field and potential patterns.



c / Two-dimensional field and potential patterns. Top: A uniformly charged rod. Bottom: A dipole. In each case, the diagram on the left shows the field vectors and constant-potential curves, while the one on the right shows the potential (up-down coordinate) as a function of x and y . Interpreting the field diagrams: Each arrow represents the field at the point where its tail has been positioned. For clarity, some of the arrows in regions of very strong field strength are not shown — they would be too long to show. Interpreting the constant-potential curves: In regions of very strong fields, the curves are not shown because they would merge together to make solid black regions. Interpreting the perspective plots: Keep in mind that even though we're visualizing things in three dimensions, these are really two-dimensional potential patterns being represented. The third (up-down) dimension represents potential, not position.

10.3 Fields by superposition

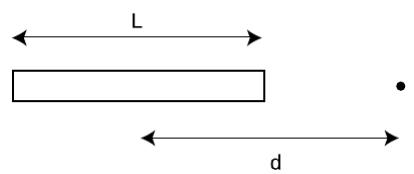
10.3.1 Electric field of a continuous charge distribution

Charge really comes in discrete chunks, but often it is mathematically convenient to treat a set of charges as if they were like a continuous fluid spread throughout a region of space. For example, a charged metal ball will have charge spread nearly uniformly all over its surface, and for most purposes it will make sense to ignore the fact that this uniformity is broken at the atomic level. The electric field made by such a continuous charge distribution is the sum of the fields created by every part of it. If we let the “parts” become infinitesimally small, we have a sum of an infinitely many infinitesimal numbers: an integral. If it was a discrete sum, as in example 3 on page 586, we would have a total electric field in the x direction that was the sum of all the x components of the individual fields, and similarly we’d have sums for the y and z components. In the continuous case, we have three integrals. Let’s keep it simple by starting with a one-dimensional example.

Field of a uniformly charged rod

example 12

- ▷ A rod of length L has charge Q spread uniformly along it. Find the electric field at a point a distance d from the center of the rod,



a / Example 12.

along the rod's axis.

▷ This is a one-dimensional situation, so we really only need to do a single integral representing the total field along the axis. We imagine breaking the rod down into short pieces of length dz , each with charge dq . Since charge is uniformly spread along the rod, we have $dq = \lambda dz$, where $\lambda = Q/L$ (Greek lambda) is the charge per unit length, in units of coulombs per meter. Since the pieces are infinitesimally short, we can treat them as point charges and use the expression $k dq/r^2$ for their contributions to the field, where $r = d - z$ is the distance from the charge at z to the point in which we are interested.

$$\begin{aligned} E_z &= \int \frac{k dq}{r^2} \\ &= \int_{-L/2}^{+L/2} \frac{k\lambda dz}{r^2} \\ &= k\lambda \int_{-L/2}^{+L/2} \frac{dz}{(d-z)^2} \end{aligned}$$

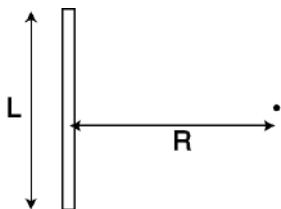
The integral can be looked up in a table, or reduced to an elementary form by substituting a new variable for $d - z$. The result is

$$\begin{aligned} E_z &= k\lambda \left(\frac{1}{d-z} \right) \Big|_{-L/2}^{+L/2} \\ &= \frac{kQ}{L} \left(\frac{1}{d-L/2} - \frac{1}{d+L/2} \right). \end{aligned}$$

For large values of d , this expression gets smaller for two reasons: (1) the denominators of the fractions become large, and (2) the two fractions become nearly the same, and tend to cancel out. This makes sense, since the field should get weaker as we get farther away from the charge. In fact, the field at large distances must approach kQ/d^2 (homework problem 2).

It's also interesting to note that the field becomes infinite at the ends of the rod, but is not infinite on the interior of the rod. Can you explain physically why this happens?

Example 12 was one-dimensional. In the general three-dimensional case, we might have to integrate all three components of the field. However, there is a trick that lets us avoid this much complication. The potential is a scalar, so we can find the potential by doing just a single integral, then use the potential to find the field.



b / Example 13.

Potential, then field

example 13

▷ A rod of length L is uniformly charged with charge Q . Find the field at a point lying in the midplane of the rod at a distance R .

▷ By symmetry, the field has only a radial component, E_R , pointing directly away from the rod (or toward it for $Q < 0$). The

brute-force approach, then, would be to evaluate the integral $E = \int |d\mathbf{E}| \cos \theta$, where $d\mathbf{E}$ is the contribution to the field from a charge dq at some point along the rod, and θ is the angle $d\mathbf{E}$ makes with the radial line.

It's easier, however, to find the potential first, and then find the field from the potential. Since the potential is a scalar, we simply integrate the contribution dV from each charge dq , without even worrying about angles and directions. Let z be the coordinate that measures distance up and down along the rod, with $z = 0$ at the center of the rod. Then the distance between a point z on the rod and the point of interest is $r = \sqrt{z^2 + R^2}$, and we have

$$\begin{aligned} V &= \int \frac{k dq}{r} \\ &= k\lambda \int_{-L/2}^{+L/2} \frac{dz}{r} \\ &= k\lambda \int_{-L/2}^{+L/2} \frac{dz}{\sqrt{z^2 + R^2}} \end{aligned}$$

The integral can be looked up in a table, or evaluated using computer software:

$$\begin{aligned} V &= k\lambda \ln \left(z + \sqrt{z^2 + R^2} \right) \Big|_{-L/2}^{+L/2} \\ &= k\lambda \ln \left(\frac{L/2 + \sqrt{L^2/4 + R^2}}{-L/2 + \sqrt{L^2/4 + R^2}} \right) \end{aligned}$$

The expression inside the parentheses can be simplified a little. Leaving out some tedious algebra, the result is

$$V = 2k\lambda \ln \left(\frac{L}{2R} + \sqrt{1 + \frac{L^2}{4R^2}} \right)$$

This can readily be differentiated to find the field:

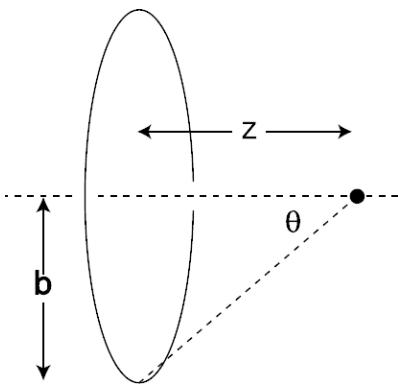
$$\begin{aligned} E_R &= -\frac{dV}{dR} \\ &= (-2k\lambda) \frac{-L/2R^2 + (1/2)(1 + L^2/4R^2)^{-1/2}(-L^2/2R^3)}{L/2R + (1 + L^2/4R^2)^{1/2}}, \end{aligned}$$

or, after some simplification,

$$E_R = \frac{k\lambda L}{R^2 \sqrt{1 + L^2/4R^2}}$$

For large values of R , the square root approaches one, and we have simply $E_R \approx k\lambda L/R^2 = kQ/R^2$. In other words, the field very far away is the same regardless of whether the charge is a point charge or some other shape like a rod. This is intuitively appealing, and doing this kind of check also helps to reassure one that the final result is correct.

The preceding example, although it involved some messy algebra, required only straightforward calculus, and no vector operations at all, because we only had to integrate a scalar function to find the potential. The next example is one in which we can integrate either the field or the potential without too much complication.



c / Example 14.

On-axis field of a ring of charge example 14

- ▷ Find the potential and field along the axis of a uniformly charged ring.
- ▷ Integrating the potential is straightforward.

$$\begin{aligned} V &= \int \frac{k dq}{r} \\ &= k \int \frac{dq}{\sqrt{b^2 + z^2}} \\ &= \frac{k}{\sqrt{b^2 + z^2}} \int dq \\ &= \frac{kQ}{\sqrt{b^2 + z^2}}, \end{aligned}$$

where Q is the total charge of the ring. This result could have been derived without calculus, since the distance r is the same for every point around the ring, i.e., the integrand is a constant. It would also be straightforward to find the field by differentiating this expression with respect to z (homework problem 10).

Instead, let's see how to find the field by direct integration. By symmetry, the field at the point of interest can have only a component along the axis of symmetry, the z axis:

$$E_x = 0$$

$$E_y = 0$$

To find the field in the z direction, we integrate the z components contributed to the field by each infinitesimal part of the ring.

$$\begin{aligned} E_z &= \int dE_z \\ &= \int |\mathbf{dE}| \cos \theta, \end{aligned}$$

where θ is the angle shown in the figure.

$$\begin{aligned} E_z &= \int \frac{k dq}{r^2} \cos \theta \\ &= k \int \frac{dq}{b^2 + z^2} \cos \theta \end{aligned}$$

Everything inside the integral is a constant, so we have

$$\begin{aligned} E_z &= \frac{k}{b^2 + z^2} \cos \theta \int dq \\ &= \frac{kQ}{b^2 + z^2} \cos \theta \\ &= \frac{kQ}{b^2 + z^2} \frac{z}{r} \\ &= \frac{kQz}{(b^2 + z^2)^{3/2}} \end{aligned}$$

In all the examples presented so far, the charge has been confined to a one-dimensional line or curve. Although it is possible, for example, to put charge on a piece of wire, it is more common to encounter practical devices in which the charge is distributed over a two-dimensional surface, as in the flat metal plates used in Thomson's experiments. Mathematically, we can approach this type of calculation with the divide-and-conquer technique: slice the surface into lines or curves whose fields we know how to calculate, and then add up the contributions to the field from all these slices. In the limit where the slices are imagined to be infinitesimally thin, we have an integral.

Field of a uniformly charged disk

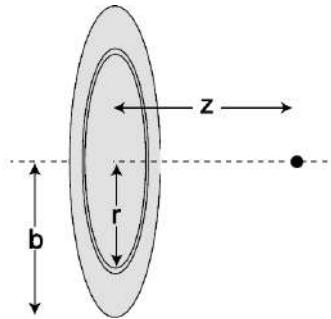
example 15

▷ A circular disk is uniformly charged. (The disk must be an insulator; if it was a conductor, then the repulsion of all the charge would cause it to collect more densely near the edge.) Find the field at a point on the axis, at a distance z from the plane of the disk.

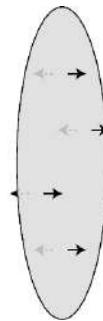
▷ We're given that every part of the disk has the same charge per unit area, so rather than working with Q , the total charge, it will be easier to use the charge per unit area, conventionally notated σ (Greek sigma), $\sigma = Q/\pi b^2$.

Since we already know the field due to a ring of charge, we can solve the problem by slicing the disk into rings, with each ring extending from r to $r + dr$. The area of such a ring equals its circumference multiplied by its width, i.e., $2\pi r dr$, so its charge is $dq = 2\pi\sigma r dr$, and from the result of example 14, its contribution to the field is

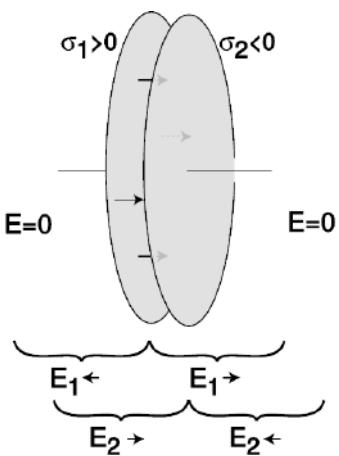
$$\begin{aligned} dE_z &= \frac{kz dq}{(r^2 + z^2)^{3/2}} \\ &= \frac{2\pi\sigma k z r dr}{(r^2 + z^2)^{3/2}} \end{aligned}$$



d / Example 15: geometry.



e / Example 15: the field on both sides (for $\sigma > 0$).



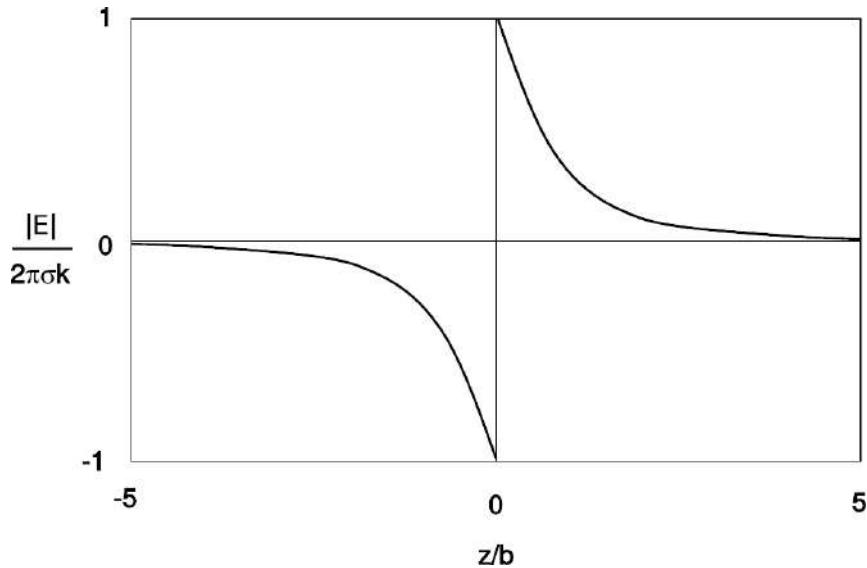
f / A capacitor consisting of two disks with opposite charges.

The total field is

$$\begin{aligned}
 E_z &= \int dE_z \\
 &= 2\pi\sigma kz \int_0^b \frac{r dr}{(r^2 + z^2)^{3/2}} \\
 &= 2\pi\sigma kz \left. \frac{-1}{\sqrt{r^2 + z^2}} \right|_{r=0}^{r=b} \\
 &= 2\pi\sigma k \left(1 - \frac{z}{\sqrt{b^2 + z^2}} \right)
 \end{aligned}$$

The result of example 15 has some interesting properties. First, we note that it was derived on the unspoken assumption of $z > 0$. By symmetry, the field on the other side of the disk must be equally strong, but in the opposite direction, as shown in figures e and g. Thus there is a discontinuity in the field at $z = 0$. In reality, the disk will have some finite thickness, and the switching over of the field will be rapid, but not discontinuous.

At large values of z , i.e., $z \gg b$, the field rapidly approaches the $1/r^2$ variation that we expect when we are so far from the disk that the disk's size and shape cannot matter (homework problem 2).



g / Example 15: variation of the field ($\sigma > 0$).

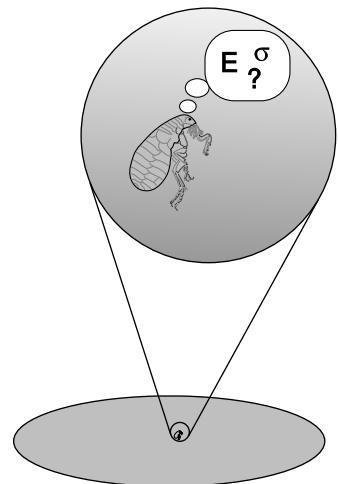
A practical application is the case of a capacitor, f, having two parallel circular plates very close together. In normal operation, the charges on the plates are opposite, so one plate has fields pointing into it and the other one has fields pointing out. In a real capacitor,

the plates are a metal conductor, not an insulator, so the charge will tend to arrange itself more densely near the edges, rather than spreading itself uniformly on each plate. Furthermore, we have only calculated the *on-axis* field in example 15; in the off-axis region, each disk's contribution to the field will be weaker, and it will also point away from the axis a little. But if we are willing to ignore these complications for the sake of a rough analysis, then the fields superimpose as shown in figure f: the fields cancel the outside of the capacitor, but between the plates its value is double that contributed by a single plate. This cancellation on the outside is a very useful property for a practical capacitor. For instance, if you look at the printed circuit board in a typical piece of consumer electronics, there are many capacitors, often placed fairly close together. If their exterior fields didn't cancel out nicely, then each capacitor would interact with its neighbors in a complicated way, and the behavior of the circuit would depend on the exact physical layout, since the interaction would be stronger or weaker depending on distance. In reality, a capacitor does create weak external electric fields, but their effects are often negligible, and we can then use the *lumped-circuit approximation*, which states that each component's behavior depends only on the currents that flow in and out of it, not on the interaction of its fields with the other components.

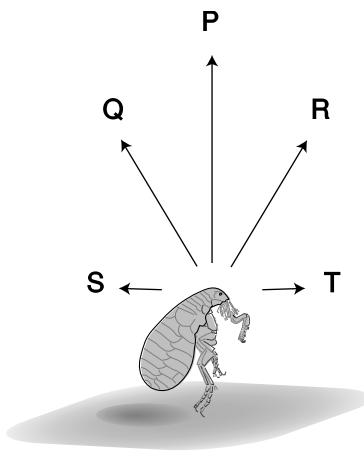
10.3.2 The field near a charged surface

From a theoretical point of view, there is something even more intriguing about example 15: the magnitude of the field for small values of z ($z \ll b$) is $E = 2\pi k\sigma$, which doesn't depend on b at all for a fixed value of σ . If we made a disk with twice the radius, and covered it with the same number of coulombs per square meter (resulting in a total charge four times as great), the field close to the disk would be unchanged! That is, a flea living near the center of the disk, h , would have no way of determining the size of her flat "planet" by measuring the local field and charge density. (Only by leaping off the surface into outer space would she be able to measure fields that were dependent on b . If she traveled very far, to $z \gg b$, she would be in the region where the field is well approximated by $|\mathbf{E}| \approx kQ/z^2 = k\pi b^2 \sigma / z^2$, which she could solve for b .)

What is the reason for this surprisingly simple behavior of the field? Is it a piece of mathematical trivia, true only in this particular case? What if the shape was a square rather than a circle? In other words, the flea gets no information about the *size* of the disk from measuring E , since $E = 2\pi k\sigma$, independent of b , but what if she didn't know the *shape*, either? If the result for a square had some other geometrical factor in front instead of 2π , then she could tell which shape it was by measuring E . The surprising mathematical fact, however, is that the result for a square, indeed for any shape whatsoever, is $E = 2\pi\sigma k$. It doesn't even matter whether the sur-



h / Close to the surface, the relationship between E and σ is a fixed one, regardless of the geometry. The flea can't determine the size or shape of her world by comparing E and σ .



i / Fields contributed by nearby parts of the surface, P, Q, and R, contribute to E_{\perp} . Fields due to distant charges, S, and T, have very small contributions to E_{\perp} because of their shallow angles.

face is flat or warped, or whether the density of charge is different at parts of the surface which are far away compared to the flea's distance above the surface.

This universal $E_{\perp} = 2\pi k\sigma$ field perpendicular to a charged surface can be proved mathematically based on Gauss's law¹ (section 10.6), but we can understand what's happening on qualitative grounds. Suppose one night, while the flea is asleep, someone adds more surface area, also positively charged, around the outside edge of her disk-shaped world, doubling its radius. The added charge, however, has very little effect on the field in her environment, as long as she stays at low altitudes above the surface. As shown in figure i, the new charge to her west contributes a field, T, that is almost purely "horizontal" (i.e., parallel to the surface) and to the east. It has a negligible upward component, since the angle is so shallow. This new eastward contribution to the field is exactly canceled out by the westward field, S, created by the new charge to her east. There is likewise almost perfect cancellation between any other pair of opposite compass directions.

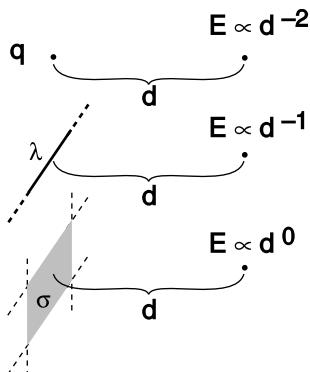
A similar argument can be made as to the shape-independence of the result, as long as the shape is symmetric. For example, suppose that the next night, the tricky real estate developers decide to add corners to the disk and transform it into a square. Each corner's contribution to the field measured at the center is canceled by the field due to the corner diagonally across from it.

What if the flea goes on a trip away from the center of the disk? The perfect cancellation of the "horizontal" fields contributed by distant charges will no longer occur, but the "vertical" field (i.e., the field perpendicular to the surface) will still be $E_{\perp} = 2\pi k\sigma$, where σ is the local charge density, since the distant charges can't contribute to the vertical field. The same result applies if the shape of the surface is asymmetric, and doesn't even have any well-defined geometric center: the component perpendicular to the surface is $E_{\perp} = 2\pi k\sigma$, but we may have $E_{\parallel} \neq 0$. All of the above arguments can be made more rigorous by discussing mathematical limits rather than using words like "very small." There is not much point in giving a rigorous proof here, however, since we will be able to demonstrate this fact as a corollary of Gauss' Law in section 10.6. The result is as follows:

At a point lying a distance z from a charged surface, the component of the electric field perpendicular to the surface obeys

$$\lim_{z \rightarrow 0} E_{\perp} = 2\pi k\sigma,$$

where σ is the charge per unit area. This is true regardless of the shape or size of the surface.



j / Example 16.

¹rhymes with "mouse"

The field near a point, line, or surface charge example 16

▷ Compare the variation of the electric field with distance, d , for small values of d in the case of a point charge, an infinite line of charge, and an infinite charged surface.

▷ For a point charge, we have already found $E \propto d^{-2}$ for the magnitude of the field, where we are now using d for the quantity we would ordinarily notate as r . This is true for all values of d , not just for small d — it has to be that way, because the point charge has no size, so if E behaved differently for small and large d , there would be no way to decide what d had to be small or large relative to.

For a line of charge, the result of example 13 is

$$E = \frac{k\lambda L}{d^2 \sqrt{1 + L^2/4d^2}}.$$

In the limit of $d \ll L$, the quantity inside the square root is dominated by the second term, and we have $E \propto d^{-1}$.

Finally, in the case of a charged surface, the result is simply $E = 2\pi\sigma k$, or $E \propto d^0$.

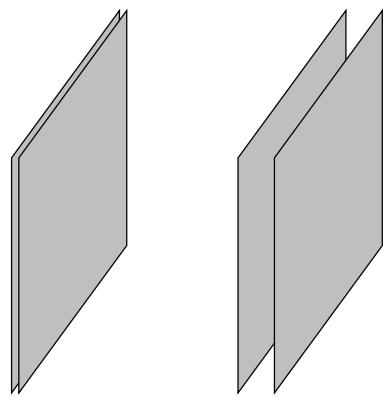
Notice the lovely simplicity of the pattern, as shown in figure j. A point is zero-dimensional: it has no length, width, or breadth. A line is one-dimensional, and a surface is two-dimensional. As the dimensionality of the charged object changes from 0 to 1, and then to 2, the exponent in the near-field expression goes from 2 to 1 to 0.

10.4 Energy in fields

10.4.1 Electric field energy

Fields possess energy, as argued on page 581, but how much energy? The answer can be found using the following elegant approach. We assume that the electric energy contained in an infinitesimal volume of space dv is given by $dU_e = f(\mathbf{E}) dv$, where f is some function, which we wish to determine, of the field \mathbf{E} . It might seem that we would have no easy way to determine the function f , but many of the functions we could cook up would violate the symmetry of space. For instance, we could imagine $f(\mathbf{E}) = aE_y$, where a is some constant with the appropriate units. However, this would violate the symmetry of space, because it would give the y axis a different status from x and z . As discussed on page 216, if we wish to calculate a scalar based on some vectors, the dot product is the only way to do it that has the correct symmetry properties. If all we have is one vector, \mathbf{E} , then the only scalar we can form is $\mathbf{E} \cdot \mathbf{E}$, which is the square of the magnitude of the electric field vector.

In principle, the energy function we are seeking could be proportional to $\mathbf{E} \cdot \mathbf{E}$, or to any function computed from it, such as $\sqrt{\mathbf{E} \cdot \mathbf{E}}$ or $(\mathbf{E} \cdot \mathbf{E})^7$. On physical grounds, however, the only possibility that works is $\mathbf{E} \cdot \mathbf{E}$. Suppose, for instance, that we pull apart two oppositely charged capacitor plates, as shown in figure a. We are doing work by pulling them apart against the force of their electrical attraction, and this quantity of mechanical work equals the increase in electrical energy, U_e . Using our previous approach to energy, we would have thought of U_e as a quantity which depended on the distance of the positive and negative charges from each other, but now we're going to imagine U_e as being stored within the electric field that exists in the space between and around the charges. When the plates are touching, their fields cancel everywhere, and there is zero electrical energy. When they are separated, there is still approximately zero field on the outside, but the field between the plates is nonzero, and holds some energy.



a / Two oppositely charged capacitor plates are pulled apart.

Now suppose we carry out the whole process, but with the plates carrying double their previous charges. Since Coulomb's law involves the product $q_1 q_2$ of two charges, we have quadrupled the force between any given pair of charged particles, and the total attractive force is therefore also four times greater than before. This means that the work done in separating the plates is four times greater, and so is the energy U_e stored in the field. The field, however, has merely been doubled at any given location: the electric field \mathbf{E}_+ due to the positively charged plate is doubled, and similarly for the contribution \mathbf{E}_- from the negative one, so the total electric field $\mathbf{E}_+ + \mathbf{E}_-$ is also doubled. Thus doubling the field results in an electrical energy which is four times greater, i.e., the energy density must be proportional to the square of the field, $dU_e \propto (\mathbf{E} \cdot \mathbf{E}) dv$. For ease

of notation, we write this as $dU_e \propto E^2 dv$, or $dU_e = aE^2 dv$, where a is a constant of proportionality. Note that we never really made use of any of the details of the geometry of figure a, so the reasoning is of general validity. In other words, not only is $dU_e = aE^2 dv$ the function that works in this particular case, but there is every reason to believe that it would work in other cases as well.

It now remains only to find a . Since the constant must be the same in all situations, we only need to find one example in which we can compute the field and the energy, and then we can determine a . The situation shown in figure a is just about the easiest example to analyze. We let the square capacitor plates be uniformly covered with charge densities $+\sigma$ and $-\sigma$, and we write b for the lengths of their sides. Let h be the gap between the plates after they have been separated. We choose $h \ll b$, so that the field experienced by the negative plate due to the positive plate is $E_+ = 2\pi k\sigma$. The charge of the negative plate is $-\sigma b^2$, so the magnitude of the force attracting it back toward the positive plate is (force) = (charge)(field) = $2\pi k\sigma^2 b^2$. The amount of work done in separating the plates is (work) = (force)(distance) = $2\pi k\sigma^2 b^2 h$. This is the amount of energy that has been stored in the field between the two plates, $U_e = 2\pi k\sigma^2 b^2 h = 2\pi k\sigma^2 v$, where v is the volume of the region between the plates.

We want to equate this to $U_e = aE^2 v$. (We can write U_e and v rather than dU_e and dv , since the field is constant in the region between the plates.) The field between the plates has contributions from both plates, $E = E_+ + E_- = 4\pi k\sigma$. (We only used half this value in the computation of the work done on the moving plate, since the moving plate can't make a force on itself. Mathematically, each plate is in a region where its own field is reversing directions, so we can think of its own contribution to the field as being zero within itself.) We then have $aE^2 v = a \cdot 16\pi^2 k^2 \sigma^2 \cdot v$, and setting this equal to $U_e = 2\pi k\sigma^2 v$ from the result of the work computation, we find $a = 1/8\pi k$. Our final result is as follows:

The electric energy possessed by an electric field \mathbf{E} occupying an infinitesimal volume of space dv is given by

$$dU_e = \frac{1}{8\pi k} E^2 dv,$$

where $E^2 = \mathbf{E} \cdot \mathbf{E}$ is the square of the magnitude of the electric field.

This is reminiscent of how waves behave: the energy content of a wave is typically proportional to the square of its amplitude.

self-check D

We can think of the quantity dU_e/dv as the *energy density* due to the electric field, i.e., the number of joules per cubic meter needed in order to create that field. (a) How does this quantity depend on the components of the field vector, E_x , E_y , and E_z ? (b) Suppose we have a field with $E_x \neq 0$, $E_y=0$, and $E_z=0$. What would happen to the energy density if we reversed the sign of E_x ?

▷ Answer, p. 1063

A numerical example

example 17

▷ A capacitor has plates whose areas are 10^{-4} m^2 , separated by a gap of 10^{-5} m . A 1.5-volt battery is connected across it. How much energy is sucked out of the battery and stored in the electric field between the plates? (A real capacitor typically has an insulating material between the plates whose molecules interact electrically with the charge in the plates. For this example, we'll assume that there is just a vacuum in between the plates. The plates are also typically rolled up rather than flat.)

▷ To connect this with our previous calculations, we need to find the charge density on the plates in terms of the voltage we were given. Our previous examples were based on the assumption that the gap between the plates was small compared to the size of the plates. Is this valid here? Well, if the plates were square, then the area of 10^{-4} m^2 would imply that their sides were 10^{-2} m in length. This is indeed very large compared to the gap of 10^{-5} m , so this assumption appears to be valid (unless, perhaps, the plates have some very strange, long and skinny shape).

Based on this assumption, the field is relatively uniform in the whole volume between the plates, so we can use a single symbol, E , to represent its magnitude, and the relation $E = dV/dx$ is equivalent to $E = \Delta V/\Delta x = (1.5 \text{ V})/(\text{gap}) = 1.5 \times 10^5 \text{ V/m}$.

Since the field is uniform, we can dispense with the calculus, and replace $dU_e = (1/8\pi k)E^2 dv$ with $U_e = (1/8\pi k)E^2 v$. The volume equals the area multiplied by the gap, so we have

$$\begin{aligned} U_e &= (1/8\pi k)E^2(\text{area})(\text{gap}) \\ &= \frac{1}{8\pi \times 9 \times 10^9 \text{ N}\cdot\text{m}^2/\text{C}^2} (1.5 \times 10^5 \text{ V/m})^2 (10^{-4} \text{ m}^2) (10^{-5} \text{ m}) \\ &= 1 \times 10^{-10} \text{ J} \end{aligned}$$

self-check E

Show that the units in the preceding example really do work out to be joules.

▷ Answer, p. 1063

Why k is on the bottom

example 18

It may also seem strange that the constant k is in the denominator of the equation $dU_e = (1/8\pi k)E^2 dv$. The Coulomb constant k tells us how strong electric forces are, so shouldn't it be on top?

No. Consider, for instance, an alternative universe in which electric forces are twice as strong as in ours. The numerical value of k is doubled. Because k is doubled, all the electric field strengths are doubled as well, which quadruples the quantity E^2 . In the expression $E^2/8\pi k$, we've quadrupled something on top and doubled something on the bottom, which makes the energy twice as big. That makes perfect sense.

Potential energy of a pair of opposite charges *example 19*

Imagine taking two opposite charges, b , that were initially far apart and allowing them to come together under the influence of their electrical attraction.

According to our old approach, electrical energy is lost because the electric force did positive work as it brought the charges together. (This makes sense because as they come together and accelerate it is their electrical energy that is being lost and converted to kinetic energy.)

By the new method, we must ask how the energy stored in the electric field has changed. In the region indicated approximately by the shading in the figure, the superposing fields of the two charges undergo partial cancellation because they are in opposing directions. The energy in the shaded region is reduced by this effect. In the unshaded region, the fields reinforce, and the energy is increased.

It would be quite a project to do an actual numerical calculation of the energy gained and lost in the two regions (this is a case where the old method of finding energy gives greater ease of computation), but it is fairly easy to convince oneself that the energy is less when the charges are closer. This is because bringing the charges together shrinks the high-energy unshaded region and enlarges the low-energy shaded region.

A spherical capacitor

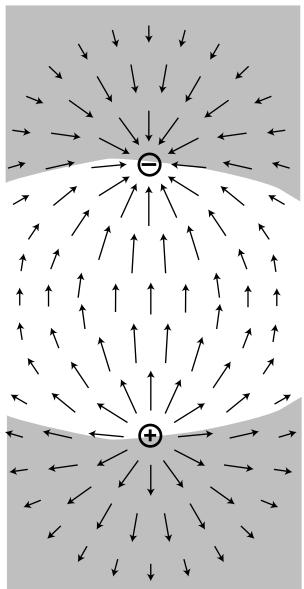
example 20

▷ A spherical capacitor, c , consists of two concentric spheres of radii a and b . Find the energy required to charge up the capacitor so that the plates hold charges $+q$ and $-q$.

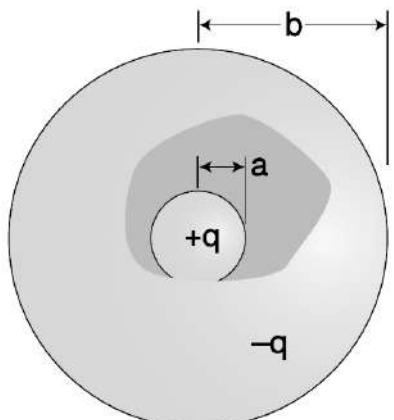
▷ On page 102, I proved that for *gravitational* forces, the interaction of a spherical shell of mass with other masses outside it is the same as if the shell's mass was concentrated at its center. On the interior of such a shell, the forces cancel out exactly. Since gravity and the electric force both vary as $1/r^2$, the same proof carries over immediately to electrical forces. The magnitude of the outward electric field contributed by the charge $+q$ of the central sphere is therefore

$$|\mathbf{E}_+| = \begin{cases} 0, & r < a \\ kq/r^2, & r > a \end{cases}$$

where r is the distance from the center. Similarly, the magnitude



b / Example 19.



c / Example B. Part of the outside sphere has been drawn as if it is transparent, in order to show the inside sphere.

of the *inward* field contributed by the outside sphere is

$$|\mathbf{E}_-| = \begin{cases} 0, & r < b \\ kq/r^2, & r > b \end{cases}.$$

In the region outside the whole capacitor, the two fields are equal in magnitude, but opposite in direction, so they cancel. We then have for the total field

$$|\mathbf{E}| = \begin{cases} 0, & r < a \\ kq/r^2, & a < r < b \\ 0, & r > b \end{cases},$$

so to calculate the energy, we only need to worry about the region $a < r < b$. The energy density in this region is

$$\begin{aligned} \frac{dU_e}{dv} &= \frac{1}{8\pi k} E^2 \\ &= \frac{kq^2}{8\pi} r^{-4}. \end{aligned}$$

This expression only depends on r , so the energy density is constant across any sphere of radius r . We can slice the region $a < r < b$ into concentric spherical layers, like an onion, and the energy within one such layer, extending from r to $r + dr$ is

$$\begin{aligned} dU_e &= \frac{dU_e}{dv} dv \\ &= \frac{dU_e}{dv} (\text{area of shell})(\text{thickness of shell}) \\ &= \left(\frac{kq^2}{8\pi} r^{-4}\right)(4\pi r^2)(dr) \\ &= \frac{kq^2}{2} r^{-2} dr. \end{aligned}$$

Integrating over all the layers to find the total energy, we have

$$\begin{aligned} U_e &= \int dU_e \\ &= \int_a^b \frac{kq^2}{2} r^{-2} dr \\ &= -\frac{kq^2}{2} r^{-1} \Big|_a^b \\ &= \frac{kq^2}{2} \left(\frac{1}{a} - \frac{1}{b}\right) \end{aligned}$$

Discussion Questions

A The figure shows a positive charge in the gap between two capacitor plates. Compare the energy of the electric fields in the two cases. Does this agree with what you would have expected based on your knowledge of electrical forces?

B The figure shows a spherical capacitor. In the text, the energy stored in its electric field is shown to be

$$U_e = \frac{kq^2}{2} \left(\frac{1}{a} - \frac{1}{b} \right).$$

What happens if the difference between b and a is very small? Does this make sense in terms of the mechanical work needed in order to separate the charges? Does it make sense in terms of the energy stored in the electric field? Should these two energies be added together?

Similarly, discuss the cases of $b \rightarrow \infty$ and $a \rightarrow 0$.

C Criticize the following statement: “A solenoid makes a charge in the space surrounding it, which dissipates when you release the energy.”

D In example 19 on page 609, I argued that for the charges shown in the figure, the fields contain less energy when the charges are closer together, because the region of cancellation expanded, while the region of reinforcing fields shrank. Perhaps a simpler approach is to consider the two extreme possibilities: the case where the charges are infinitely far apart, and the one in which they are at zero distance from each other, i.e., right on top of each other. Carry out this reasoning for the case of (1) a positive charge and a negative charge of equal magnitude, (2) two positive charges of equal magnitude, (3) the gravitational energy of two equal masses.

10.4.2 Gravitational field energy

Example B depended on the close analogy between electric and gravitational forces. In fact, every argument, proof, and example discussed so far in this section is equally valid as a gravitational example, provided we take into account one fact: only positive mass exists, and the gravitational force between two masses is attractive. This is the opposite of what happens with electrical forces, which are repulsive in the case of two positive charges. As a consequence of this, we need to assign a *negative* energy density to the gravitational field! For a gravitational field, we have

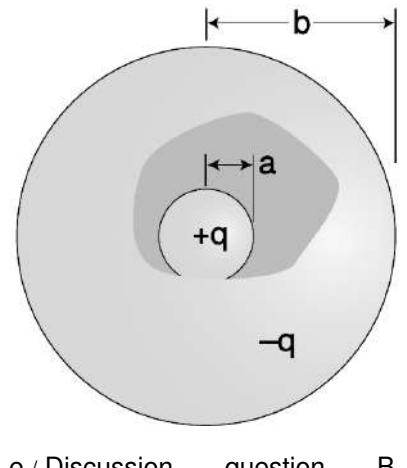
$$dU_g = -\frac{1}{8\pi G} g^2 dv,$$

where $g^2 = \mathbf{g} \cdot \mathbf{g}$ is the square of the magnitude of the gravitational field.

10.4.3 Magnetic field energy

So far we've only touched in passing on the topic of magnetic fields, which will deal with in detail in chapter 11. Magnetism is an interaction between moving charge and moving charge, i.e., between currents and currents. Since a current has a direction in

$\begin{array}{c} + + + + + + \\ \hline \cdot + \end{array}$ \hline (1)	$\begin{array}{c} + + + + + + \\ \hline \cdot + \end{array}$ \hline (2)	
d / Discussion	question	A.



e / Discussion question B.

space,² while charge doesn't, we can anticipate that the mathematical rule connecting a magnetic field to its source-currents will have to be completely different from the one relating the electric field to its source-charges. However, if you look carefully at the argument leading to the relation $dU_e/dv = E^2/8\pi k$, you'll see that these mathematical details were only necessary to the part of the argument in which we fixed the constant of proportionality. To establish $dU_e/dv \propto E^2$, we only had to use three simple facts:

- The field is proportional to the source.
- Forces are proportional to fields.
- Field contributed by multiple sources add like vectors.

All three of these statements are true for the magnetic field as well, so without knowing anything more specific about magnetic fields — not even what units are used to measure them! — we can state with certainty that the energy density in the magnetic field is proportional to the square of the magnitude of the magnetic field. The constant of proportionality is given on p. 693.

²Current is a scalar, since the definition $I = dq/dt$ is the derivative of a scalar. However, there is a closely related quantity called the current *density*, \mathbf{J} , which is a vector, and \mathbf{J} is in fact the more fundamentally important quantity.

10.5 LRC circuits

The long road leading from the light bulb to the computer started with one very important step: the introduction of feedback into electronic circuits. Although the principle of feedback has been understood and applied to mechanical systems for centuries, and to electrical ones since the early twentieth century, for most of us the word evokes an image of Jimi Hendrix intentionally creating earsplitting screeches, or of the school principal doing the same inadvertently in the auditorium. In the guitar example, the musician stands in front of the amp and turns it up so high that the sound waves coming from the speaker come back to the guitar string and make it shake harder. This is an example of *positive* feedback: the harder the string vibrates, the stronger the sound waves, and the stronger the sound waves, the harder the string vibrates. The only limit is the power-handling ability of the amplifier.

Negative feedback is equally important. Your thermostat, for example, provides negative feedback by kicking the heater off when the house gets warm enough, and by firing it up again when it gets too cold. This causes the house's temperature to oscillate back and forth within a certain range. Just as out-of-control exponential freak-outs are a characteristic behavior of positive-feedback systems, oscillation is typical in cases of negative feedback. You have already studied negative feedback extensively in section 3.3 in the case of a mechanical system, although we didn't call it that.

10.5.1 Capacitance and inductance

In a mechanical oscillation, energy is exchanged repetitively between potential and kinetic forms, and may also be siphoned off in the form of heat dissipated by friction. In an electrical circuit, resistors are the circuit elements that dissipate heat. What are the electrical analogs of storing and releasing the potential and kinetic energy of a vibrating object? When you think of energy storage in an electrical circuit, you are likely to imagine a battery, but even rechargeable batteries can only go through 10 or 100 cycles before they wear out. In addition, batteries are not able to exchange energy on a short enough time scale for most applications. The circuit in a musical synthesizer may be called upon to oscillate thousands of times a second, and your microwave oven operates at gigahertz frequencies. Instead of batteries, we generally use capacitors and inductors to store energy in oscillating circuits. Capacitors, which you've already encountered, store energy in electric fields. An inductor does the same with magnetic fields.

Capacitors

A capacitor's energy exists in its surrounding electric fields. It is proportional to the square of the field strength, which is proportional to the charges on the plates. If we assume the plates carry charges



a / The symbol for a capacitor.



b / Some capacitors.

that are the same in magnitude, $+q$ and $-q$, then the energy stored in the capacitor must be proportional to q^2 . For historical reasons, we write the constant of proportionality as $1/2C$,

$$U_C = \frac{1}{2C} q^2.$$

The constant C is a geometrical property of the capacitor, called its capacitance.

Based on this definition, the units of capacitance must be coulombs squared per joule, and this combination is more conveniently abbreviated as the farad, $1 \text{ F} = 1 \text{ C}^2/\text{J}$. “Condenser” is a less formal term for a capacitor. Note that the labels printed on capacitors often use MF to mean μF , even though MF should really be the symbol for megafarads, not microfarads. Confusion doesn’t result from this nonstandard notation, since picofarad and microfarad values are the most common, and it wasn’t until the 1990’s that even millifarad and farad values became available in practical physical sizes. Figure a shows the symbol used in schematics to represent a capacitor.

A parallel-plate capacitor

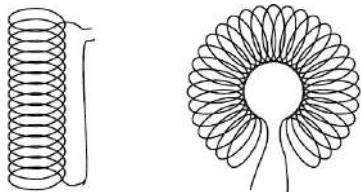
example 21

▷ Suppose a capacitor consists of two parallel metal plates with area A , and the gap between them is h . The gap is small compared to the dimensions of the plates. What is the capacitance?

▷ Since the plates are metal, the charges on each plate are free to move, and will tend to cluster themselves more densely near the edges due to the mutual repulsion of the other charges in the same plate. However, it turns out that if the gap is small, this is a small effect, so we can get away with assuming uniform charge density on each plate. The result of example 17 then applies, and for the region between the plates, we have $E = 4\pi k\sigma = 4\pi kq/A$ and $U_e = (1/8\pi k)E^2 Ah$. Substituting the first expression into the second, we find $U_e = 2\pi kq^2 h/A$. Comparing this to the definition of capacitance, we end up with $C = A/4\pi kh$.

Inductors

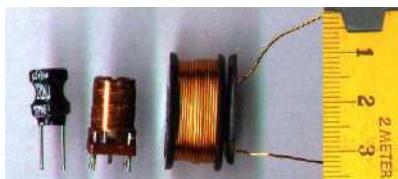
Any current will create a magnetic field, so in fact every current-carrying wire in a circuit acts as an inductor! However, this type of “stray” inductance is typically negligible, just as we can usually ignore the stray resistance of our wires and only take into account the actual resistors. To store any appreciable amount of magnetic energy, one usually uses a coil of wire designed specifically to be an inductor. All the loops’ contribution to the magnetic field add together to make a stronger field. Unlike capacitors and resistors, practical inductors are easy to make by hand. One can for instance spool some wire around a short wooden dowel. An inductor like this, in the form cylindrical coil of wire, is called a solenoid, c, and a stylized solenoid, d, is the symbol used to represent an inductor



c / Two common geometries for inductors. The cylindrical shape on the left is called a solenoid.



d / The symbol for an inductor.



e / Some inductors.

in a circuit regardless of its actual geometry.

How much energy does an inductor store? The energy density is proportional to the square of the magnetic field strength, which is in turn proportional to the current flowing through the coiled wire, so the energy stored in the inductor must be proportional to I^2 . We write $L/2$ for the constant of proportionality, giving

$$U_L = \frac{L}{2} I^2.$$

As in the definition of capacitance, we have a factor of $1/2$, which is purely a matter of definition. The quantity L is called the *inductance* of the inductor, and we see that its units must be joules per ampere squared. This clumsy combination of units is more commonly abbreviated as the henry, $1 \text{ henry} = 1 \text{ J/A}^2$. Rather than memorizing this definition, it makes more sense to derive it when needed from the definition of inductance. Many people know inductors simply as “coils,” or “chokes,” and will not understand you if you refer to an “inductor,” but they will still refer to L as the “inductance,” not the “coilance” or “chokeance!”

There is a lumped circuit approximation for inductors, just like the one for capacitors (p. 603). For a capacitor, this means assuming that the electric fields are completely internal, so that components only interact via currents that flow through wires, not due to the physical overlapping of their fields in space. Similarly for an inductor, the lumped circuit approximation is the assumption that the magnetic fields are completely internal.

Identical inductances in series

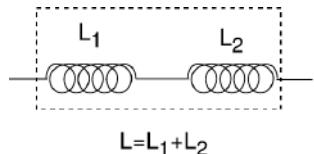
example 22

If two inductors are placed in series, any current that passes through the combined double inductor must pass through both its parts. If we assume the lumped circuit approximation, the two inductors’ fields don’t interfere with each other, so the energy is doubled for a given current. Thus by the definition of inductance, the inductance is doubled as well. In general, inductances in series add, just like resistances. The same kind of reasoning also shows that the inductance of a solenoid is approximately proportional to its length, assuming the number of turns per unit length is kept constant. (This is only approximately true, because putting two solenoids end-to-end causes the fields just outside their mouths to overlap and add together in a complicated manner. In other words, the lumped-circuit approximation may not be very good.)

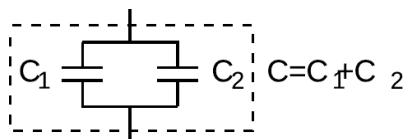
Identical capacitances in parallel

example 23

When two identical capacitances are placed in parallel, any charge deposited at the terminals of the combined double capacitor will divide itself evenly between the two parts. The electric fields surrounding each capacitor will be half the intensity, and therefore



f / Inductances in series add.



g / Capacitances in parallel add.

store one quarter the energy. Two capacitors, each storing one quarter the energy, give half the total energy storage. Since capacitance is inversely related to energy storage, this implies that identical capacitances in parallel give double the capacitance. In general, capacitances in parallel add. This is unlike the behavior of inductors and resistors, for which series configurations give addition.

This is consistent with the result of example 21, which had the capacitance of a single parallel-plate capacitor proportional to the area of the plates. If we have two parallel-plate capacitors, and we combine them in parallel and bring them very close together side by side, we have produced a single capacitor with plates of double the area, and it has approximately double the capacitance, subject to any violation of the lumped-circuit approximation due to the interaction of the fields where the edges of the capacitors are joined together.

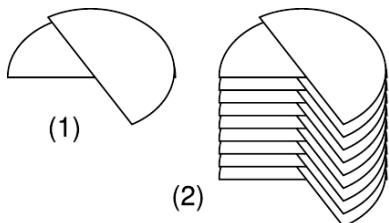
Inductances in parallel and capacitances in series are explored in homework problems 36 and 33.

A variable capacitor

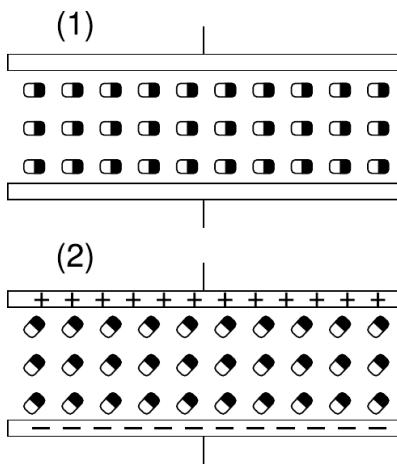
example 24

Figure h/1 shows the construction of a variable capacitor out of two parallel semicircles of metal. One plate is fixed, while the other can be rotated about their common axis with a knob. The opposite charges on the two plates are attracted to one another, and therefore tend to gather in the overlapping area. This overlapping area, then, is the only area that effectively contributes to the capacitance, and turning the knob changes the capacitance. The simple design can only provide very small capacitance values, so in practice one usually uses a bank of capacitors, wired in parallel, with all the moving parts on the same shaft.

Discussion Questions



h / A variable capacitor.



i / Discussion question B.

A Suppose that two parallel-plate capacitors are wired in parallel, and are placed very close together, side by side, so that the lumped circuit approximation is not very accurate. Will the resulting capacitance be too small, or too big? Could you twist the circuit into a different shape and make the effect be the other way around, or make the effect vanish? How about the case of two inductors in series?

B Most practical capacitors do not have an air gap or vacuum gap between the plates; instead, they have an insulating substance called a dielectric. We can think of the molecules in this substance as dipoles that are free to rotate (at least a little), but that are not free to move around, since it is a solid. The figure shows a highly stylized and unrealistic way of visualizing this. We imagine that all the dipoles are initially turned sideways, (1), and that as the capacitor is charged, they all respond by turning through a certain angle, (2). (In reality, the scene might be much more random, and the alignment effect much weaker.)

For simplicity, imagine inserting just one electric dipole into the vacuum gap. For a given amount of charge on the plates, how does this affect

the amount of energy stored in the electric field? How does this affect the capacitance?

Now redo the analysis in terms of the mechanical work needed in order to charge up the plates.

10.5.2 Oscillations

Figure j shows the simplest possible oscillating circuit. For any useful application it would actually need to include more components. For example, if it was a radio tuner, it would need to be connected to an antenna and an amplifier. Nevertheless, all the essential physics is there.

We can analyze it without any sweat or tears whatsoever, simply by constructing an analogy with a mechanical system. In a mechanical oscillator, k , we have two forms of stored energy,

$$U_{\text{spring}} = \frac{1}{2}kx^2 \quad (1)$$

$$K = \frac{1}{2}mv^2. \quad (2)$$

In the case of a mechanical oscillator, we have usually assumed a friction force of the form that turns out to give the nicest mathematical results, $F = -bv$. In the circuit, the dissipation of energy into heat occurs via the resistor, with no mechanical force involved, so in order to make the analogy, we need to restate the role of the friction force in terms of energy. The power dissipated by friction equals the mechanical work it does in a time interval dt , divided by dt , $P = W/dt = F dx/dt = Fv = -bv^2$, so

$$\text{rate of heat dissipation} = -bv^2. \quad (3)$$

self-check F

Equation (1) has x squared, and equations (2) and (3) have v squared. Because they're squared, the results don't depend on whether these variables are positive or negative. Does this make physical sense? ▷ Answer, p. 1063

In the circuit, the stored forms of energy are

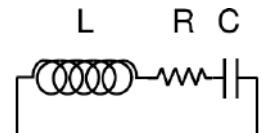
$$U_C = \frac{1}{2C}q^2 \quad (1')$$

$$U_L = \frac{1}{2}LI^2, \quad (2')$$

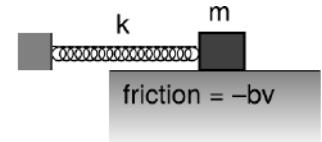
and the rate of heat dissipation in the resistor is

$$\text{rate of heat dissipation} = -RI^2. \quad (3')$$

Comparing the two sets of equations, we first form analogies between quantities that represent the state of the system at some moment



j / A series LRC circuit.



k / A mechanical analogy for the LRC circuit.

in time:

$$\begin{aligned}x &\leftrightarrow q \\v &\leftrightarrow I\end{aligned}$$

self-check G

How is v related mathematically to x ? How is I connected to q ? Are the two relationships analogous? ▷ Answer, p. 1063

Next we relate the ones that describe the system's permanent characteristics:

$$\begin{aligned}k &\leftrightarrow 1/C \\m &\leftrightarrow L \\b &\leftrightarrow R\end{aligned}$$

Since the mechanical system naturally oscillates with a frequency³ $\omega \approx \sqrt{k/m}$, we can immediately solve the electrical version by analogy, giving

$$\omega \approx \frac{1}{\sqrt{LC}}.$$

Since the resistance R is analogous to b in the mechanical case, we find that the Q (quality factor, not charge) of the resonance is inversely proportional to R , and the width of the resonance is directly proportional to R .

Tuning a radio receiver

example 25

A radio receiver uses this kind of circuit to pick out the desired station. Since the receiver resonates at a particular frequency, stations whose frequencies are far off will not excite any response in the circuit. The value of R has to be small enough so that only one station at a time is picked up, but big enough so that the tuner isn't too touchy. The resonant frequency can be tuned by adjusting either L or C , but variable capacitors are easier to build than variable inductors.

A numerical calculation

example 26

The phone company sends more than one conversation at a time over the same wire, which is accomplished by shifting each voice signal into different range of frequencies during transmission. The number of signals per wire can be maximized by making each range of frequencies (known as a bandwidth) as small as possible. It turns out that only a relatively narrow range of frequencies is necessary in order to make a human voice intelligible, so the phone company filters out all the extreme highs and lows. (This is why your phone voice sounds different from your normal voice.)

³As in chapter 2, we use the word "frequency" to mean either f or $\omega = 2\pi f$ when the context makes it clear which is being referred to.

- ▷ If the filter consists of an LRC circuit with a broad resonance centered around 1.0 kHz, and the capacitor is 1 μF (microfarad), what inductance value must be used?

▷ Solving for L , we have

$$\begin{aligned} L &= \frac{1}{C\omega^2} \\ &= \frac{1}{(10^{-6} \text{ F})(2\pi \times 10^3 \text{ s}^{-1})^2} \\ &= 2.5 \times 10^{-3} \text{ F}^{-1}\text{s}^2 \end{aligned}$$

Checking that these really are the same units as henries is a little tedious, but it builds character:

$$\begin{aligned} \text{F}^{-1}\text{s}^2 &= (\text{C}^2/\text{J})^{-1}\text{s}^2 \\ &= \text{J} \cdot \text{C}^{-2}\text{s}^2 \\ &= \text{J}/\text{A}^2 \\ &= \text{H} \end{aligned}$$

The result is 25 mH (millihenries).

This is actually quite a large inductance value, and would require a big, heavy, expensive coil. In fact, there is a trick for making this kind of circuit small and cheap. There is a kind of silicon chip called an op-amp, which, among other things, can be used to simulate the behavior of an inductor. The main limitation of the op-amp is that it is restricted to low-power applications.

10.5.3 Voltage and current

What is physically happening in one of these oscillating circuits? Let's first look at the mechanical case, and then draw the analogy to the circuit. For simplicity, let's ignore the existence of damping, so there is no friction in the mechanical oscillator, and no resistance in the electrical one.

Suppose we take the mechanical oscillator and pull the mass away from equilibrium, then release it. Since friction tends to resist the spring's force, we might naively expect that having zero friction would allow the mass to leap instantaneously to the equilibrium position. This can't happen, however, because the mass would have to have infinite velocity in order to make such an instantaneous leap. Infinite velocity would require infinite kinetic energy, but the only kind of energy that is available for conversion to kinetic is the energy stored in the spring, and that is finite, not infinite. At each step on its way back to equilibrium, the mass's velocity is controlled exactly by the amount of the spring's energy that has so far been converted into kinetic energy. After the mass reaches equilibrium, it overshoots due to its own momentum. It performs identical oscillations on both sides of equilibrium, and it never loses amplitude because friction is not available to convert mechanical energy into heat.

Now with the electrical oscillator, the analog of position is charge. Pulling the mass away from equilibrium is like depositing charges $+q$ and $-q$ on the plates of the capacitor. Since resistance tends to resist the flow of charge, we might imagine that with no friction present, the charge would instantly flow through the inductor (which is, after all, just a piece of wire), and the capacitor would discharge instantly. However, such an instant discharge is impossible, because it would require infinite current for one instant. Infinite current would create infinite magnetic fields surrounding the inductor, and these fields would have infinite energy. Instead, the rate of flow of current is controlled at each instant by the relationship between the amount of energy stored in the magnetic field and the amount of current that must exist in order to have that strong a field. After the capacitor reaches $q = 0$, it overshoots. The circuit has its own kind of electrical “inertia,” because if charge was to stop flowing, there would have to be zero current through the inductor. But the current in the inductor must be related to the amount of energy stored in its magnetic fields. When the capacitor is at $q = 0$, all the circuit’s energy is in the inductor, so it must therefore have strong magnetic fields surrounding it and quite a bit of current going through it.

The only thing that might seem spooky here is that we used to speak as if the current in the inductor caused the magnetic field, but now it sounds as if the field causes the current. Actually this is symptomatic of the elusive nature of cause and effect in physics. It’s equally valid to think of the cause and effect relationship in either way. This may seem unsatisfying, however, and for example does not really get at the question of what brings about a voltage difference across the resistor (in the case where the resistance is finite); there must be such a voltage difference, because without one, Ohm’s law would predict zero current through the resistor.

Voltage, then, is what is really missing from our story so far.

Let’s start by studying the voltage across a capacitor. Voltage is electrical potential energy per unit charge, so the voltage difference between the two plates of the capacitor is related to the amount by which its energy would increase if we increased the absolute values of the charges on the plates from q to $q + dq$:

$$\begin{aligned} V_C &= (U_{q+dq} - U_q) / dq \\ &= \frac{dU_C}{dq} \\ &= \frac{d}{dq} \left(\frac{1}{2C} q^2 \right) \\ &= \frac{q}{C} \end{aligned}$$

Many books use this as the definition of capacitance. This equation, by the way, probably explains the historical reason why C was de-

fined so that the energy was *inversely* proportional to C for a given value of q : the people who invented the definition were thinking of a capacitor as a device for storing charge rather than energy, and the amount of charge stored for a fixed voltage (the charge “capacity”) is proportional to C .

In the case of an inductor, we know that if there is a steady, constant current flowing through it, then the magnetic field is constant, and so is the amount of energy stored; no energy is being exchanged between the inductor and any other circuit element. But what if the current is changing? The magnetic field is proportional to the current, so a change in one implies a change in the other. For concreteness, let’s imagine that the magnetic field and the current are both decreasing. The energy stored in the magnetic field is therefore decreasing, and by conservation of energy, this energy can’t just go away — some other circuit element must be taking energy from the inductor. The simplest example, shown in figure 1, is a series circuit consisting of the inductor plus one other circuit element. It doesn’t matter what this other circuit element is, so we just call it a black box, but if you like, we can think of it as a resistor, in which case the energy lost by the inductor is being turned into heat by the resistor. The junction rule tells us that both circuit elements have the same current through them, so I could refer to either one, and likewise the loop rule tells us $V_{\text{inductor}} + V_{\text{black box}} = 0$, so the two voltage drops have the same absolute value, which we can refer to as V . Whatever the black box is, the rate at which it is taking energy from the inductor is given by $|P| = |IV|$, so

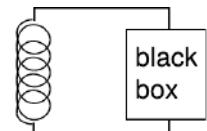
$$\begin{aligned} |IV| &= \left| \frac{dU_L}{dt} \right| \\ &= \left| \frac{d}{dt} \left(\frac{1}{2} LI^2 \right) \right| \\ &= \left| LI \frac{dI}{dt} \right|, \end{aligned}$$

or

$$|V| = \left| L \frac{dI}{dt} \right|,$$

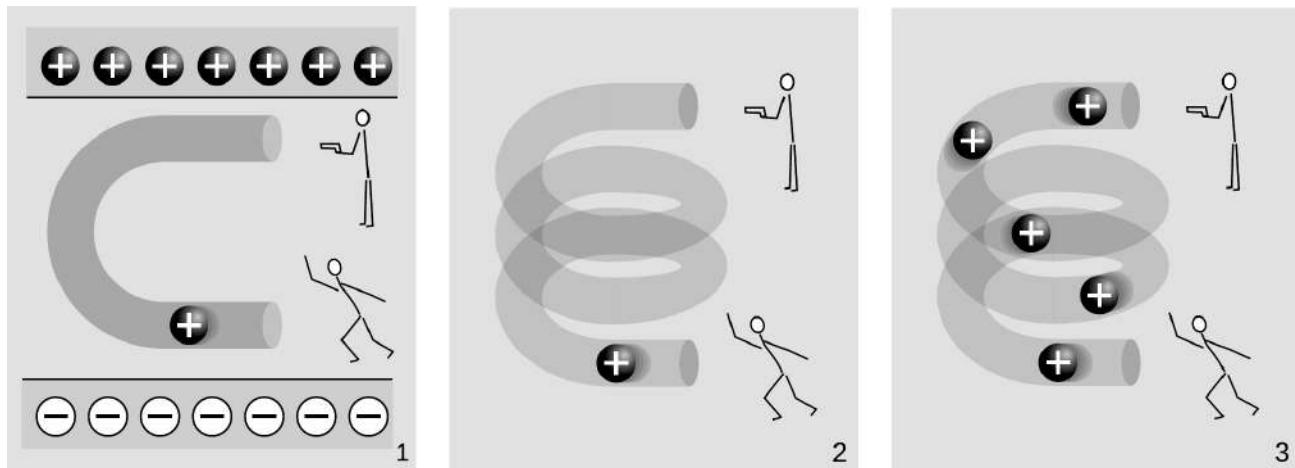
which in many books is taken to be the definition of inductance. The direction of the voltage drop (plus or minus sign) is such that the inductor resists the change in current.

There’s one very intriguing thing about this result. Suppose, for concreteness, that the black box in figure 1 is a resistor, and that the inductor’s energy is decreasing, and being converted into heat in the resistor. The voltage drop across the resistor indicates that it has an electric field across it, which is driving the current.



I / The inductor releases energy and gives it to the black box.

But where is this electric field coming from? There are no charges anywhere that could be creating it! What we've discovered is one special case of a more general principle, the principle of induction: a changing magnetic field creates an electric field, which is in addition to any electric field created by charges. (The reverse is also true: any electric field that changes over time creates a magnetic field.) Induction forms the basis for such technologies as the generator and the transformer, and ultimately it leads to the existence of light, which is a wave pattern in the electric and magnetic fields. These are all topics for chapter 11, but it's truly remarkable that we could come to this conclusion without yet having learned any details about magnetism.



m / Electric fields made by charges, 1, and by changing magnetic fields, 2 and 3.

The cartoons in figure m compares electric fields made by charges, 1, to electric fields made by changing magnetic fields, 2-3. In m/1, two physicists are in a room whose ceiling is positively charged and whose floor is negatively charged. The physicist on the bottom throws a positively charged bowling ball into the curved pipe. The physicist at the top uses a radar gun to measure the speed of the ball as it comes out of the pipe. They find that the ball has slowed down by the time it gets to the top. By measuring the change in the ball's kinetic energy, the two physicists are acting just like a voltmeter. They conclude that the top of the tube is at a higher voltage than the bottom of the pipe. A difference in voltage indicates an electric field, and this field is clearly being caused by the charges in the floor and ceiling.

In m/2, there are no charges anywhere in the room except for the charged bowling ball. Moving charges make magnetic fields, so there is a magnetic field surrounding the helical pipe while the ball is moving through it. A magnetic field has been created where there

was none before, and that field has energy. Where could the energy have come from? It can only have come from the ball itself, so the ball must be losing kinetic energy. The two physicists working together are again acting as a voltmeter, and again they conclude that there is a voltage difference between the top and bottom of the pipe. This indicates an electric field, but this electric field can't have been created by any charges, because there aren't any in the room. This electric field was created by the change in the magnetic field.

The bottom physicist keeps on throwing balls into the pipe, until the pipe is full of balls, $m/3$, and finally a steady current is established. While the pipe was filling up with balls, the energy in the magnetic field was steadily increasing, and that energy was being stolen from the balls' kinetic energy. But once a steady current is established, the energy in the magnetic field is no longer changing. The balls no longer have to give up energy in order to build up the field, and the physicist at the top finds that the balls are exiting the pipe at full speed again. There is no voltage difference any more. Although there is a current, dI/dt is zero.

Ballasts

example 27

In a gas discharge tube, such as a neon sign, enough voltage is applied to a tube full of gas to ionize some of the atoms in the gas. Once ions have been created, the voltage accelerates them, and they strike other atoms, ionizing them as well and resulting in a chain reaction. This is a spark, like a bolt of lightning. But once the spark starts up, the device begins to act as though it has no resistance: more and more current flows, without the need to apply any more voltage. The power, $P = IV$, would grow without limit, and the tube would burn itself out.

The simplest solution is to connect an inductor, known as the "ballast," in series with the tube, and run the whole thing on an AC voltage. During each cycle, as the voltage reaches the point where the chain reaction begins, there is a surge of current, but the inductor resists such a sudden change of current, and the energy that would otherwise have burned out the bulb is instead channeled into building a magnetic field.

A common household fluorescent lightbulb consists of a gas discharge tube in which the glass is coated with a fluorescent material. The gas in the tube emits ultraviolet light, which is absorbed by the coating, and the coating then glows in the visible spectrum.

Until recently, it was common for a fluorescent light's ballast to be a simple inductor, and for the whole device to be operated at the 60 Hz frequency of the electrical power lines. This caused the lights to flicker annoyingly at 120 Hz, and could also cause an audible hum, since the magnetic field surrounding the inductor could



n / Ballasts for fluorescent lights. Top: a big, heavy inductor used as a ballast in an old-fashioned fluorescent bulb. Bottom: a small solid-state ballast, built into the base of a modern compact fluorescent bulb.

exert mechanical forces on things. Modern compact fluorescent bulbs have ballasts built into their bases that use a frequency in the kilohertz range, eliminating the flicker and hum.

Discussion Question

- A** What happens when the physicist at the bottom in figure m/3 starts getting tired, and decreases the current?

10.5.4 Decay

Up until now I've soft-pedaled the fact that by changing the characteristics of an oscillator, it is possible to produce non-oscillatory behavior. For example, imagine taking the mass-on-a-spring system and making the spring weaker and weaker. In the limit of small k , it's as though there was no spring whatsoever, and the behavior of the system is that if you kick the mass, it simply starts slowing down. For friction proportional to v , as we've been assuming, the result is that the velocity approaches zero, but never actually reaches zero. This is unrealistic for the mechanical oscillator, which will not have vanishing friction at low velocities, but it is quite realistic in the case of an electrical circuit, for which the voltage drop across the resistor really does approach zero as the current approaches zero.

We do not even have to reduce k to exactly zero in order to get non-oscillatory behavior. There is actually a finite, critical value below which the behavior changes, so that the mass never even makes it through one cycle. This is the case of overdamping, discussed on page 190.

Electrical circuits can exhibit all the same behavior. For simplicity we will analyze only the cases where either the capacitor or the inductor is completely absent, giving $Q = 0$.

The RC circuit

We first analyze the RC circuit, o. In reality one would have to "kick" the circuit, for example by briefly inserting a battery, in order to get any interesting behavior. We start with Ohm's law and the equation for the voltage across a capacitor:

$$V_R = IR$$

$$V_C = q/C$$

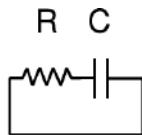
The loop rule tells us

$$V_R + V_C = 0,$$

and combining the three equations results in a relationship between q and I :

$$I = -\frac{1}{RC}q$$

The negative sign tells us that the current tends to reduce the charge on the capacitor, i.e., to discharge it. It makes sense that the current



o / An RC circuit.

is proportional to q : if q is large, then the attractive forces between the $+q$ and $-q$ charges on the plates of the capacitor are large, and charges will flow more quickly through the resistor in order to reunite. If there was zero charge on the capacitor plates, there would be no reason for current to flow. Since amperes, the unit of current, are the same as coulombs per second, it appears that the quantity RC must have units of seconds, and you can check for yourself that this is correct. RC is therefore referred to as the time constant of the circuit.

How exactly do I and q vary with time? Rewriting I as dq/dt , we have

$$\frac{dq}{dt} = -\frac{1}{RC}q.$$

We need a function $q(t)$ whose derivative equals itself, but multiplied by a negative constant. A function of the form ae^t , where $e = 2.718\dots$ is the base of natural logarithms, is the only one that has its derivative equal to itself, and ae^{bt} has its derivative equal to itself multiplied by b . Thus our solution is

$$q = q_0 \exp\left(-\frac{t}{RC}\right).$$

The RL circuit

The RL circuit, q , can be attacked by similar methods, and it can easily be shown that it gives

$$I = I_0 \exp\left(-\frac{R}{L}t\right).$$

The RL time constant equals L/R .

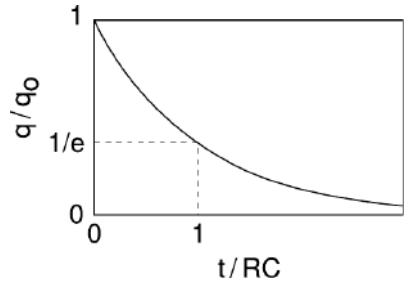
Death by solenoid; spark plugs

example 28

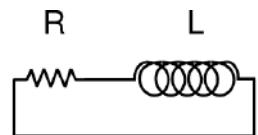
When we suddenly break an RL circuit, what will happen? It might seem that we're faced with a paradox, since we only have two forms of energy, magnetic energy and heat, and if the current stops suddenly, the magnetic field must collapse suddenly. But where does the lost magnetic energy go? It can't go into resistive heating of the resistor, because the circuit has now been broken, and current can't flow!

The way out of this conundrum is to recognize that the open gap in the circuit has a resistance which is large, but not infinite. This large resistance causes the RL time constant L/R to be very small. The current thus continues to flow for a very brief time, and flows straight across the air gap where the circuit has been opened. In other words, there is a spark!

We can determine based on several different lines of reasoning that the voltage drop from one end of the spark to the other must



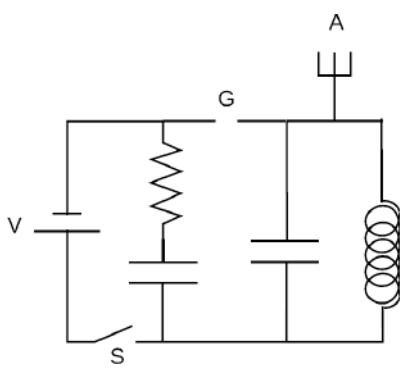
p / Over a time interval RC , the charge on the capacitor is reduced by a factor of e .



q / An RL circuit.

be very large. First, the air's resistance is large, so $V = IR$ requires a large voltage. We can also reason that all the energy in the magnetic field is being dissipated in a short time, so the power dissipated in the spark, $P = IV$, is large, and this requires a large value of V . (I isn't large — it is decreasing from its initial value.) Yet a third way to reach the same result is to consider the equation $V_L = dI/dt$: since the time constant is short, the time derivative dI/dt is large.

This is exactly how a car's spark plugs work. Another application is to electrical safety: it can be dangerous to break an inductive circuit suddenly, because so much energy is released in a short time. There is also no guarantee that the spark will discharge across the air gap; it might go through your body instead, since your body might have a lower resistance.



r / Example 29.

A spark-gap radio transmitter

example 29

Figure r shows a primitive type of radio transmitter, called a spark gap transmitter, used to send Morse code around the turn of the twentieth century. The high voltage source, V , is typically about 10,000 volts. When the telegraph switch, S , is closed, the RC circuit on the left starts charging up. An increasing voltage difference develops between the electrodes of the spark gap, G . When this voltage difference gets large enough, the electric field in the air between the electrodes causes a spark, partially discharging the RC circuit, but charging the LC circuit on the right. The LC circuit then oscillates at its resonant frequency (typically about 1 MHz), but the energy of these oscillations is rapidly radiated away by the antenna, A , which sends out radio waves (chapter 11).

Discussion Questions

- A** A gopher gnaws through one of the wires in the DC lighting system in your front yard, and the lights turn off. At the instant when the circuit becomes open, we can consider the bare ends of the wire to be like the plates of a capacitor, with an air gap (or gopher gap) between them. What kind of capacitance value are we talking about here? What would this tell you about the RC time constant?

10.5.5 Review of complex numbers

For a more detailed treatment of complex numbers, see ch. 3 of James Nearing's free book at physics.miami.edu/~nearing/mathmethods.

We assume there is a number, i , such that $i^2 = -1$. The square roots of -1 are then i and $-i$. (In electrical engineering work, where i stands for current, j is sometimes used instead.) This gives rise to a number system, called the complex numbers, containing the real numbers as a subset. Any complex number z can be written in the form $z = a + bi$, where a and b are real, and a and b are then referred to as the real and imaginary parts of z . A number with a zero real part is called an imaginary number. The complex numbers can be visualized as a plane, with the real number line placed horizontally like the x axis of the familiar $x-y$ plane, and the imaginary numbers running along the y axis. The complex numbers are complete in a way that the real numbers aren't: every nonzero complex number has two square roots. For example, 1 is a real number, so it is also a member of the complex numbers, and its square roots are -1 and 1 . Likewise, -1 has square roots i and $-i$, and the number i has square roots $1/\sqrt{2} + i/\sqrt{2}$ and $-1/\sqrt{2} - i/\sqrt{2}$.

Complex numbers can be added and subtracted by adding or subtracting their real and imaginary parts. Geometrically, this is the same as vector addition.

The complex numbers $a + bi$ and $a - bi$, lying at equal distances above and below the real axis, are called complex conjugates. The results of the quadratic formula are either both real, or complex conjugates of each other. The complex conjugate of a number z is denoted as \bar{z} or z^* .

The complex numbers obey all the same rules of arithmetic as the reals, except that they can't be ordered along a single line. That is, it's not possible to say whether one complex number is greater than another. We can compare them in terms of their magnitudes (their distances from the origin), but two distinct complex numbers may have the same magnitude, so, for example, we can't say whether 1 is greater than i or i is greater than 1 .

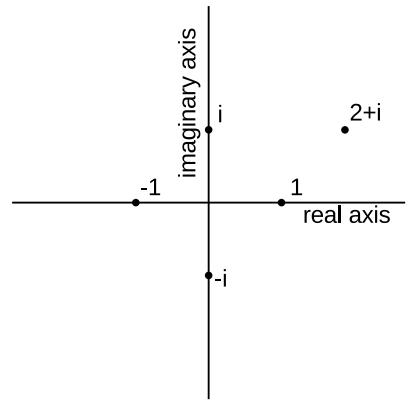
A square root of i

example 30

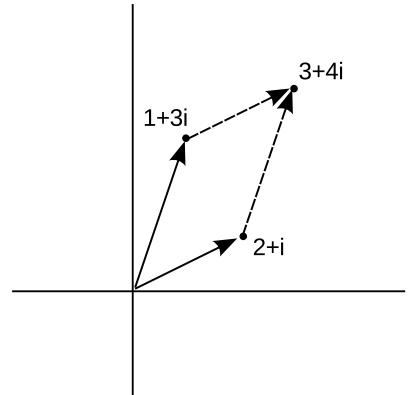
▷ Prove that $1/\sqrt{2} + i/\sqrt{2}$ is a square root of i .

▷ Our proof can use any ordinary rules of arithmetic, except for ordering.

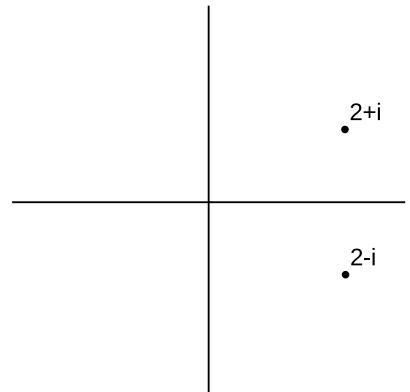
$$\begin{aligned} \left(\frac{1}{\sqrt{2}} + \frac{i}{\sqrt{2}}\right)^2 &= \frac{1}{\sqrt{2}} \cdot \frac{1}{\sqrt{2}} + \frac{1}{\sqrt{2}} \cdot \frac{i}{\sqrt{2}} + \frac{i}{\sqrt{2}} \cdot \frac{1}{\sqrt{2}} + \frac{i}{\sqrt{2}} \cdot \frac{i}{\sqrt{2}} \\ &= \frac{1}{2}(1 + i + i - 1) \\ &= i \end{aligned}$$



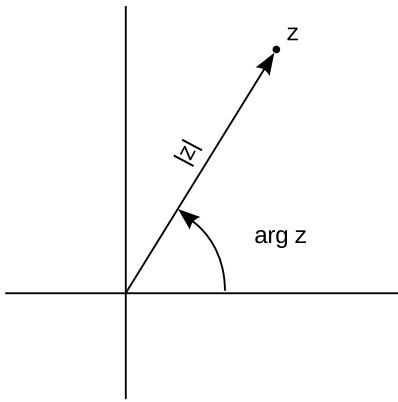
s / Visualizing complex numbers as points in a plane.



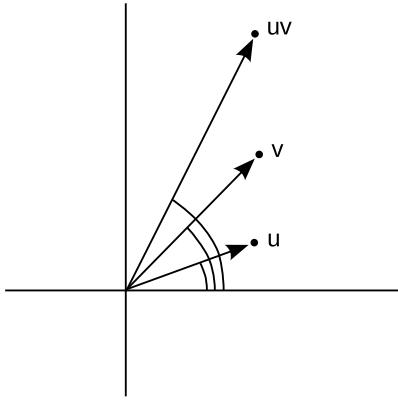
t / Addition of complex numbers is just like addition of vectors, although the real and imaginary axes don't actually represent directions in space.



u / A complex number and its conjugate.



w / A complex number can be described in terms of its magnitude and argument.



w / The argument of uv is the sum of the arguments of u and v .

Example 30 showed one method of multiplying complex numbers. However, there is another nice interpretation of complex multiplication. We define the argument of a complex number as its angle in the complex plane, measured counterclockwise from the positive real axis. Multiplying two complex numbers then corresponds to multiplying their magnitudes, and adding their arguments.

self-check H

Using this interpretation of multiplication, how could you find the square roots of a complex number? ▷ Answer, p. 1063

An identity

example 31

The magnitude $|z|$ of a complex number z obeys the identity $|z|^2 = z\bar{z}$. To prove this, we first note that \bar{z} has the same magnitude as z , since flipping it to the other side of the real axis doesn't change its distance from the origin. Multiplying z by \bar{z} gives a result whose magnitude is found by multiplying their magnitudes, so the magnitude of $z\bar{z}$ must therefore equal $|z|^2$. Now we just have to prove that $z\bar{z}$ is a positive real number. But if, for example, z lies counterclockwise from the real axis, then \bar{z} lies clockwise from it. If z has a positive argument, then \bar{z} has a negative one, or vice-versa. The sum of their arguments is therefore zero, so the result has an argument of zero, and is on the positive real axis.
4

This whole system was built up in order to make every number have square roots. What about cube roots, fourth roots, and so on? Does it get even more weird when you want to do those as well? No. The complex number system we've already discussed is sufficient to handle all of them. The nicest way of thinking about it is in terms of roots of polynomials. In the real number system, the polynomial $x^2 - 1$ has two roots, i.e., two values of x (plus and minus one) that we can plug in to the polynomial and get zero. Because it has these two real roots, we can rewrite the polynomial as $(x - 1)(x + 1)$. However, the polynomial $x^2 + 1$ has no real roots. It's ugly that in the real number system, some second-order polynomials have two roots, and can be factored, while others can't. In the complex number system, they all can. For instance, $x^2 + 1$ has roots i and $-i$, and can be factored as $(x - i)(x + i)$. In general, the fundamental theorem of algebra states that in the complex number system, any n th-order polynomial can be factored completely into n linear factors, and we can also say that it has n complex roots, with the understanding that some of the roots may be the same. For instance, the fourth-order polynomial $x^4 + x^2$ can be factored as $(x - i)(x + i)(x - 0)(x - 0)$, and we say that it has four roots, i , $-i$, 0 , and 0 , two of which happen to be the same. This is a sensible way to think about it, because

⁴I cheated a little. If z 's argument is 30 degrees, then we could say \bar{z} 's was -30, but we could also call it 330. That's OK, because $330 + 30$ gives 360, and an argument of 360 is the same as an argument of zero.

in real life, numbers are always approximations anyway, and if we make tiny, random changes to the coefficients of this polynomial, it will have four distinct roots, of which two just happen to be very close to zero.

Discussion Questions

A Find $\arg i$, $\arg(-i)$, and $\arg 37$, where $\arg z$ denotes the argument of the complex number z .

B Visualize the following multiplications in the complex plane using the interpretation of multiplication in terms of multiplying magnitudes and adding arguments: $(i)(i) = -1$, $(i)(-i) = 1$, $(-i)(-i) = -1$.

C If we visualize z as a point in the complex plane, how should we visualize $-z$? What does this mean in terms of arguments? Give similar interpretations for z^2 and \sqrt{z} .

D Find four different complex numbers z such that $z^4 = 1$.

E Compute the following. For the final two, use the magnitude and argument, not the real and imaginary parts.

$$|1+i| , \quad \arg(1+i) , \quad \left| \frac{1}{1+i} \right| , \quad \arg\left(\frac{1}{1+i}\right) ,$$

From these, find the real and imaginary parts of $1/(1+i)$.

10.5.6 Euler's formula

Having expanded our horizons to include the complex numbers, it's natural to want to extend functions we knew and loved from the world of real numbers so that they can also operate on complex numbers. The only really natural way to do this in general is to use Taylor series. A particularly beautiful thing happens with the functions e^x , $\sin x$, and $\cos x$:

$$\begin{aligned} e^x &= 1 + \frac{1}{2!}x^2 + \frac{1}{3!}x^3 + \dots \\ \cos x &= 1 - \frac{1}{2!}x^2 + \frac{1}{4!}x^4 - \dots \\ \sin x &= x - \frac{1}{3!}x^3 + \frac{1}{5!}x^5 - \dots \end{aligned}$$

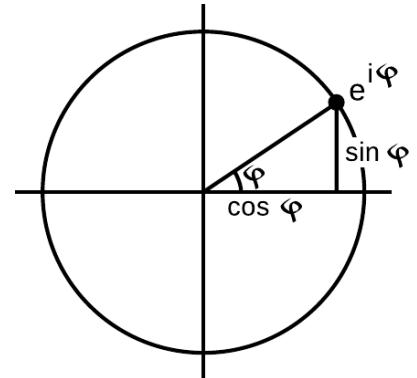
If $x = i\phi$ is an imaginary number, we have

$$e^{i\phi} = \cos \phi + i \sin \phi,$$

a result known as Euler's formula. The geometrical interpretation in the complex plane is shown in figure x.

Although the result may seem like something out of a freak show at first, applying the definition of the exponential function makes it clear how natural it is:

$$e^x = \lim_{n \rightarrow \infty} \left(1 + \frac{x}{n}\right)^n.$$



x / The complex number $e^{i\phi}$ lies on the unit circle.



y / Leonhard Euler (1707-1783)

When $x = i\phi$ is imaginary, the quantity $(1 + i\phi/n)$ represents a number lying just above 1 in the complex plane. For large n , $(1 + i\phi/n)$ becomes very close to the unit circle, and its argument is the small angle ϕ/n . Raising this number to the n th power multiplies its argument by n , giving a number with an argument of ϕ .

Euler's formula is used frequently in physics and engineering.

Trig functions in terms of complex exponentials *example 32*

- ▷ Write the sine and cosine functions in terms of exponentials.
- ▷ Euler's formula for $x = -i\phi$ gives $\cos \phi - i \sin \phi$, since $\cos(-\theta) = \cos \theta$, and $\sin(-\theta) = -\sin \theta$.

$$\cos x = \frac{e^{ix} + e^{-ix}}{2}$$

$$\sin x = \frac{e^{ix} - e^{-ix}}{2i}$$

A hard integral made easy *example 33*

- ▷ Evaluate

$$\int e^x \cos x \, dx$$

- ▷ This seemingly impossible integral becomes easy if we rewrite the cosine in terms of exponentials:

$$\begin{aligned} & \int e^x \cos x \, dx \\ &= \int e^x \left(\frac{e^{ix} + e^{-ix}}{2} \right) \, dx \\ &= \frac{1}{2} \int (e^{(1+i)x} + e^{(1-i)x}) \, dx \\ &= \frac{1}{2} \left(\frac{e^{(1+i)x}}{1+i} + \frac{e^{(1-i)x}}{1-i} \right) + C \end{aligned}$$

Since this result is the integral of a real-valued function, we'd like it to be real, and in fact it is, since the first and second terms are complex conjugates of one another. If we wanted to, we could use Euler's theorem to convert it back to a manifestly real result.⁵

10.5.7 Impedance

So far we have been thinking in terms of the free oscillations of a circuit. This is like a mechanical oscillator that has been kicked but then left to oscillate on its own without any external force to keep the vibrations from dying out. Suppose an LRC circuit is driven with a sinusoidally varying voltage, such as will occur when a radio

⁵In general, the use of complex number techniques to do an integral could result in a complex number, but that complex number would be a constant, which could be subsumed within the usual constant of integration.

tuner is hooked up to a receiving antenna. We know that a current will flow in the circuit, and we know that there will be resonant behavior, but it is not necessarily simple to relate current to voltage in the most general case. Let's start instead with the special cases of LRC circuits consisting of only a resistance, only a capacitance, or only an inductance. We are interested only in the steady-state response.

The purely resistive case is easy. Ohm's law gives

$$I = \frac{V}{R}.$$

In the purely capacitive case, the relation $V = q/C$ lets us calculate

$$\begin{aligned} I &= \frac{dq}{dt} \\ &= C \frac{dV}{dt}. \end{aligned}$$

This is partly analogous to Ohm's law. For example, if we double the amplitude of a sinusoidally varying AC voltage, the derivative dV/dt will also double, and the amplitude of the sinusoidally varying current will also double. However, it is not true that $I = V/R$, because taking the derivative of a sinusoidal function shifts its phase by 90 degrees. If the voltage varies as, for example, $V(t) = V_0 \sin(\omega t)$, then the current will be $I(t) = \omega C V_0 \cos(\omega t)$. The amplitude of the current is $\omega C V_0$, which is proportional to V_0 , but it's not true that $I(t) = V(t)/R$ for some constant R .

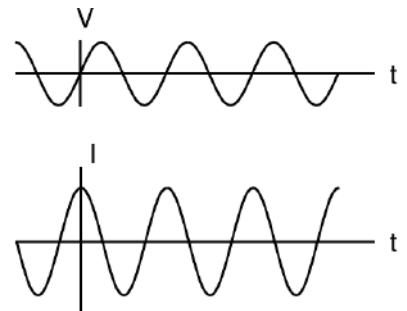
A second problem that crops up is that our entire analysis of DC resistive circuits was built on the foundation of the loop rule and the junction rule, both of which are statements about sums. To apply the junction rule to an AC circuit, for example, we would say that the sum of the sine waves describing the currents coming into the junction is equal (at every moment in time) to the sum of the sine waves going out. Now sinusoidal functions have a remarkable property, which is that if you add two different sinusoidal functions having the same frequency, the result is also a sinusoid with that frequency. For example, $\cos \omega t + \sin \omega t = \sqrt{2} \sin(\omega t + \pi/4)$, which can be proved using trig identities. The trig identities can get very cumbersome, however, and there is a much easier technique involving complex numbers.

Figure aa shows a useful way to visualize what's going on. When a circuit is oscillating at a frequency ω , we use points in the plane to represent sinusoidal functions with various phases and amplitudes.

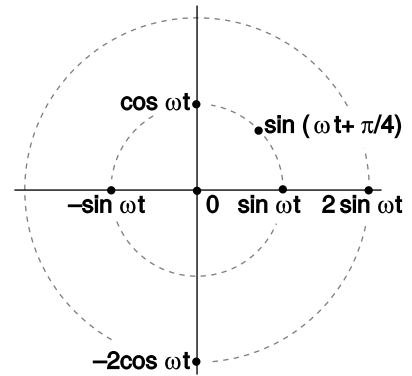
self-check 1

Which of the following functions can be represented in this way? $\cos(6t - 4)$, $\cos^2 t$, $\tan t$

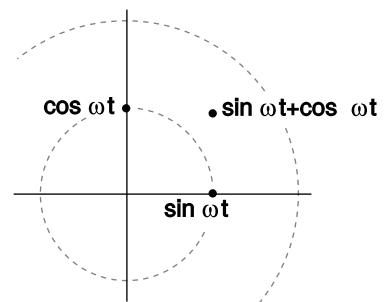
▷ Answer, p. 1064



z / In a capacitor, the current is 90° ahead of the voltage in phase.



aa / Representing functions with points in polar coordinates.



ab / Adding two sinusoidal functions.

The simplest examples of how to visualize this in polar coordinates are ones like $\cos \omega t + \cos \omega t = 2 \cos \omega t$, where everything has the same phase, so all the points lie along a single line in the polar plot, and addition is just like adding numbers on the number line. The less trivial example $\cos \omega t + \sin \omega t = \sqrt{2} \sin(\omega t + \pi/4)$, can be visualized as in figure ab.

Figure ab suggests that all of this can be tied together nicely if we identify our plane with the plane of complex numbers. For example, the complex numbers 1 and i represent the functions $\sin \omega t$ and $\cos \omega t$. In figure z, for example, the voltage across the capacitor is a sine wave multiplied by a number that gives its amplitude, so we associate that function with a number \tilde{V} lying on the real axis. Its magnitude, $|\tilde{V}|$, gives the amplitude in units of volts, while its argument $\arg \tilde{V}$, gives its phase angle, which is zero. The current is a multiple of the cosine, so we identify it with a number \tilde{I} lying on the imaginary axis. We have $\arg \tilde{I} = 90^\circ$, and $|\tilde{I}|$ is the amplitude of the current, in units of amperes. But comparing with our result above, we have $|\tilde{I}| = \omega C |\tilde{V}|$. Bringing together the phase and magnitude information, we have $\tilde{I} = i\omega C \tilde{V}$. This looks very much like Ohm's law, so we write

$$\tilde{I} = \frac{\tilde{V}}{Z_C},$$

where the quantity

$$Z_C = -\frac{i}{\omega C}, \quad [\text{impedance of a capacitor}]$$

having units of ohms, is called the *impedance* of the capacitor at this frequency.

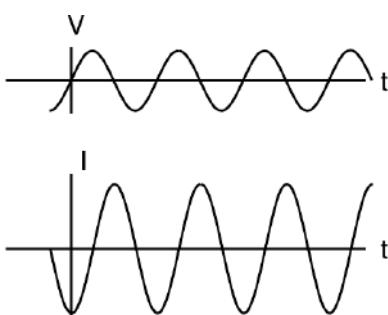
It makes sense that the impedance becomes infinite at zero frequency. Zero frequency means that it would take an infinite time before the voltage would change by any amount. In other words, this is like a situation where the capacitor has been connected across the terminals of a battery and been allowed to settle down to a state where there is constant charge on both terminals. Since the electric fields between the plates are constant, there is no energy being added to or taken out of the field. A capacitor that can't exchange energy with any other circuit component is nothing more than a broken (open) circuit.

Note that we have two types of complex numbers: those that represent sinusoidal functions of time, and those that represent impedances. The ones that represent sinusoidal functions have tildes on top, which look like little sine waves.

self-check J

Why can't a capacitor have its impedance printed on it along with its capacitance?

▷ Answer, p. 1064



ac / The current through an inductor lags behind the voltage by a phase angle of 90° .

Similar math (but this time with an integral instead of a derivative) gives

$$Z_L = i\omega L \quad [\text{impedance of an inductor}]$$

for an inductor. It makes sense that the inductor has lower impedance at lower frequencies, since at zero frequency there is no change in the magnetic field over time. No energy is added to or released from the magnetic field, so there are no induction effects, and the inductor acts just like a piece of wire with negligible resistance. The term “choke” for an inductor refers to its ability to “choke out” high frequencies.

The phase relationships shown in figures z and ac can be remembered using my own mnemonic, “eVIL,” which shows that the voltage (V) leads the current (I) in an inductive circuit, while the opposite is true in a capacitive one. A more traditional mnemonic is “ELI the ICE man,” which uses the notation E for emf, a concept closely related to voltage (see p. 715).

Summarizing, the impedances of resistors, capacitors, and inductors are

$$\begin{aligned} Z_R &= R \\ Z_C &= -\frac{i}{\omega C} \\ Z_L &= i\omega L. \end{aligned}$$

Low-pass and high-pass filters

example 34

An LRC circuit only responds to a certain range (band) of frequencies centered around its resonant frequency. As a filter, this is known as a bandpass filter. If you turn down both the bass and the treble on your stereo, you have created a bandpass filter.

To create a high-pass or low-pass filter, we only need to insert a capacitor or inductor, respectively, in series. For instance, a very basic surge protector for a computer could be constructed by inserting an inductor in series with the computer. The desired 60 Hz power from the wall is relatively low in frequency, while the surges that can damage your computer show much more rapid time variation. Even if the surges are not sinusoidal signals, we can think of a rapid “spike” qualitatively as if it was very high in frequency — like a high-frequency sine wave, it changes very rapidly.

Inductors tend to be big, heavy, expensive circuit elements, so a simple surge protector would be more likely to consist of a capacitor in *parallel* with the computer. (In fact one would normally just connect one side of the power circuit to ground via a capacitor.) The capacitor has a very high impedance at the low frequency of the desired 60 Hz signal, so it siphons off very little of the current.

But for a high-frequency signal, the capacitor's impedance is very small, and it acts like a zero-impedance, easy path into which the current is diverted.

The main things to be careful about with impedance are that (1) the concept only applies to a circuit that is being driven sinusoidally, (2) the impedance of an inductor or capacitor is frequency-dependent.

Discussion Question

A Figure z on page 631 shows the voltage and current for a capacitor. Sketch the q - t graph, and use it to give a physical explanation of the phase relationship between the voltage and current. For example, why is the current zero when the voltage is at a maximum or minimum?

B Figure ac on page 632 shows the voltage and current for an inductor. The power is considered to be positive when energy is being put into the inductor's magnetic field. Sketch the graph of the power, and then the graph of U , the energy stored in the magnetic field, and use it to give a physical explanation of the P - t graph. In particular, discuss why the frequency is doubled on the P - t graph.

C Relate the features of the graph in figure ac on page 632 to the story told in cartoons in figure m/2-3 on page 622.

10.5.8 Power

How much power is delivered when an oscillating voltage is applied to an impedance? The equation $P = IV$ is generally true, since voltage is defined as energy per unit charge, and current is defined as charge per unit time: multiplying them gives energy per unit time. In a DC circuit, all three quantities were constant, but in an oscillating (AC) circuit, all three display time variation.

A resistor

First let's examine the case of a resistor. For instance, you're probably reading this book from a piece of paper illuminated by a glowing lightbulb, which is driven by an oscillating voltage with amplitude V_0 . In the special case of a resistor, we know that I and V are in phase. For example, if V varies as $V_0 \cos \omega t$, then I will be a cosine as well, $I_0 \cos \omega t$. The power is then $I_0 V_0 \cos^2 \omega t$, which is always positive,⁶ and varies between 0 and $I_0 V_0$. Even if the time variation was $\cos \omega t$ or $\sin(\omega t + \pi/4)$, we would still have a maximum power of $I_0 V_0$, because both the voltage and the current would reach their maxima at the same time. In a lightbulb, the moment of maximum power is when the circuit is most rapidly heating the filament. At the instant when $P = 0$, a quarter of a cycle later, no current is flowing, and no electrical energy is being turned into heat. Throughout the whole cycle, the filament is getting rid of energy by

⁶A resistor always turns electrical energy into heat. It never turns heat into electrical energy!

radiating light.⁷ Since the circuit oscillates at a frequency⁸ of 60 Hz, the temperature doesn't really have time to cycle up or down very much over the 1/60 s period of the oscillation, and we don't notice any significant variation in the brightness of the light, even with a short-exposure photograph.

Thus, what we really want to know is the average power, “average” meaning the average over one full cycle. Since we’re covering a whole cycle with our average, it doesn’t matter what phase we assume. Let’s use a cosine. The total amount of energy transferred over one cycle is

$$\begin{aligned} E &= \int dE \\ &= \int_0^T \frac{dE}{dt} dt, \end{aligned}$$

where $T = 2\pi/\omega$ is the period.

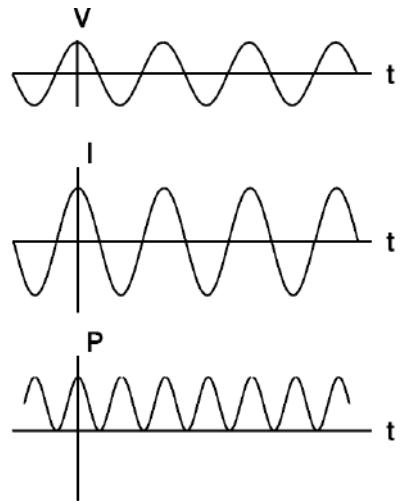
$$\begin{aligned} E &= \int_0^T P dt \\ &= \int_0^T I_o V_o \cos^2 \omega t dt \\ &= I_o V_o \int_0^T \cos^2 \omega t dt \\ &= I_o V_o \int_0^T \frac{1}{2} (1 + \cos 2\omega t) dt \end{aligned}$$

The reason for using the trig identity $\cos^2 x = (1 + \cos 2x)/2$ in the last step is that it lets us get the answer without doing a hard integral. Over the course of one full cycle, the quantity $\cos 2\omega t$ goes positive, negative, positive, and negative again, so the integral of it is zero. We then have

$$\begin{aligned} E &= I_o V_o \int_0^T \frac{1}{2} dt \\ &= \frac{I_o V_o T}{2} \end{aligned}$$

⁷To many people, the word “radiation” implies nuclear contamination. Actually, the word simply means something that “radiates” outward. Natural sunlight is “radiation.” So is the light from a lightbulb, or the infrared light being emitted by your skin right now.

⁸Note that this time “frequency” means f , not ω ! Physicists and engineers generally use ω because it simplifies the equations, but electricians and technicians always use f . The 60 Hz frequency is for the U.S.



ad / Power in a resistor: the rate at which electrical energy is being converted into heat.

The average power is

$$\begin{aligned} P_{av} &= \frac{\text{energy transferred in one full cycle}}{\text{time for one full cycle}} \\ &= \frac{I_o V_o T / 2}{T} \\ &= \frac{I_o V_o}{2}, \end{aligned}$$

i.e., the average is half the maximum. The power varies from 0 to $I_o V_o$, and it spends equal amounts of time above and below the maximum, so it isn't surprising that the average power is half-way in between zero and the maximum. Summarizing, we have

$$P_{av} = \frac{I_o V_o}{2} \quad [\text{average power in a resistor}]$$

for a resistor.

RMS quantities

Suppose one day the electric company decided to start supplying your electricity as DC rather than AC. How would the DC voltage have to be related to the amplitude V_o of the AC voltage previously used if they wanted your lightbulbs to have the same brightness as before? The resistance of the bulb, R , is a fixed value, so we need to relate the power to the voltage and the resistance, eliminating the current. In the DC case, this gives $P = IV = (V/R)V = V^2/R$. (For DC, P and P_{av} are the same.) In the AC case, $P_{av} = I_o V_o / 2 = V_o^2 / 2R$. Since there is no factor of 1/2 in the DC case, the same power could be provided with a DC voltage that was smaller by a factor of $1/\sqrt{2}$. Although you will hear people say that household voltage in the U.S. is 110 V, its amplitude is actually $(110 \text{ V}) \times \sqrt{2} \approx 160 \text{ V}$. The reason for referring to $V_o/\sqrt{2}$ as "the" voltage is that people who are naive about AC circuits can plug $V_o/\sqrt{2}$ into a familiar DC equation like $P = V^2/R$ and get the right *average* answer. The quantity $V_o/\sqrt{2}$ is called the "RMS" voltage, which stands for "root mean square." The idea is that if you square the function $V(t)$, take its average (mean) over one cycle, and then take the square root of that average, you get $V_o/\sqrt{2}$. Many digital meters provide RMS readouts for measuring AC voltages and currents.

A capacitor

For a capacitor, the calculation starts out the same, but ends up with a twist. If the voltage varies as a cosine, $V_o \cos \omega t$, then the relation $I = C dV/dt$ tells us that the current will be some constant multiplied by minus the sine, $-V_o \sin \omega t$. The integral we did in the case of a resistor now becomes

$$E = \int_0^T -I_o V_o \sin \omega t \cos \omega t dt,$$

and based on figure ae, you can easily convince yourself that over the course of one full cycle, the power spends two quarter-cycles being negative and two being positive. In other words, the average power is zero!

Why is this? It makes sense if you think in terms of energy. A resistor converts electrical energy to heat, never the other way around. A capacitor, however, merely stores electrical energy in an electric field and then gives it back. For a capacitor,

$$P_{av} = 0 \quad [\text{average power in a capacitor}]$$

Notice that although the average power is zero, the power at any given instant is *not* typically zero, as shown in figure ae. The capacitor *does* transfer energy: it's just that after borrowing some energy, it always pays it back in the next quarter-cycle.

An inductor

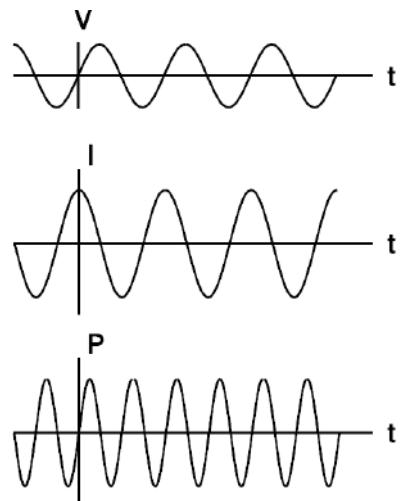
The analysis for an inductor is similar to that for a capacitor: the power averaged over one cycle is zero. Again, we're merely storing energy temporarily in a field (this time a magnetic field) and getting it back later.

10.5.9 Impedance matching

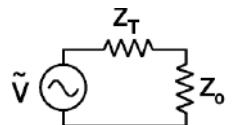
Figure af shows a commonly encountered situation: we wish to maximize the average power, P_{av} , delivered to the load for a fixed value of V_o , the amplitude of the oscillating driving voltage. We assume that the impedance of the transmission line, Z_T is a fixed value, over which we have no control, but we are able to design the load, Z_o , with any impedance we like. For now, we'll also assume that both impedances are resistive. For example, Z_T could be the resistance of a long extension cord, and Z_o could be a lamp at the end of it. The result generalizes immediately, however, to any kind of impedance. For example, the load could be a stereo speaker's magnet coil, which is displays both inductance and resistance. (For a purely inductive or capacitive load, P_{av} equals zero, so the problem isn't very interesting!)

Since we're assuming both the load and the transmission line are resistive, their impedances add in series, and the amplitude of the current is given by

$$I_o = \frac{V_o}{Z_o + Z_T},$$



ae / Power in a capacitor: the rate at which energy is being stored in (+) or removed from (-) the electric field.



af / We wish to maximize the power delivered to the load, Z_o , by adjusting its impedance.

so

$$\begin{aligned} P_{av} &= I_o V_o / 2 \\ &= I_o^2 Z_o / 2 \\ &= \frac{V_o^2 Z_o}{(Z_o + Z_T)^2} / 2. \end{aligned}$$

The maximum of this expression occurs where the derivative is zero,

$$\begin{aligned} 0 &= \frac{1}{2} \frac{d}{dZ_o} \left[\frac{V_o^2 Z_o}{(Z_o + Z_T)^2} \right] \\ 0 &= \frac{1}{2} \frac{d}{dZ_o} \left[\frac{Z_o}{(Z_o + Z_T)^2} \right] \\ 0 &= (Z_o + Z_T)^{-2} - 2Z_o (Z_o + Z_T)^{-3} \\ 0 &= (Z_o + Z_T) - 2Z_o \\ Z_o &= Z_T \end{aligned}$$

In other words, to maximize the power delivered to the load, we should make the load's impedance the same as the transmission line's. This result may seem surprising at first, but it makes sense if you think about it. If the load's impedance is too high, it's like opening a switch and breaking the circuit; no power is delivered. On the other hand, it doesn't pay to make the load's impedance too small. Making it smaller does give more current, but no matter how small we make it, the current will still be limited by the transmission line's impedance. As the load's impedance approaches zero, the current approaches this fixed value, and the the power delivered, $I_o^2 Z_o$, decreases in proportion to Z_o .

Maximizing the power transmission by matching Z_T to Z_o is called *impedance matching*. For example, an 8-ohm home stereo speaker will be correctly matched to a home stereo amplifier with an internal impedance of 8 ohms, and 4-ohm car speakers will be correctly matched to a car stereo with a 4-ohm internal impedance. You might think impedance matching would be unimportant because even if, for example, we used a car stereo to drive 8-ohm speakers, we could compensate for the mismatch simply by turning the volume knob higher. This is indeed one way to compensate for any impedance mismatch, but there is always a price to pay. When the impedances are matched, half the power is dissipated in the transmission line and half in the load. By connecting a 4-ohm amplifier to an 8-ohm speaker, however, you would be setting up a situation in two watts were being dissipated as heat inside the amp for every watt being delivered to the speaker. In other words, you would be wasting energy, and perhaps burning out your amp when you turned up the volume to compensate for the mismatch.

10.5.10 Impedances in series and parallel

How do impedances combine in series and parallel? The beauty of treating them as complex numbers is that they simply combine according to the same rules you've already learned as resistances.

Series impedance

example 35

- ▷ A capacitor and an inductor in series with each other are driven by a sinusoidally oscillating voltage. At what frequency is the current maximized?
- ▷ Impedances in series, like resistances in series, add. The capacitor and inductor act as if they were a single circuit element with an impedance

$$\begin{aligned} Z &= Z_L + Z_C \\ &= i\omega L - \frac{i}{\omega C}. \end{aligned}$$

The current is then

$$\tilde{I} = \frac{\tilde{V}}{i\omega L - i/\omega C}.$$

We don't care about the phase of the current, only its amplitude, which is represented by the absolute value of the complex number \tilde{I} , and this can be maximized by making $|i\omega L - i/\omega C|$ as small as possible. But there is some frequency at which this quantity is zero —

$$\begin{aligned} 0 &= i\omega L - \frac{i}{\omega C} \\ \frac{1}{\omega C} &= \omega L \\ \omega &= \frac{1}{\sqrt{LC}} \end{aligned}$$

At this frequency, the current is infinite! What is going on physically? This is an LRC circuit with $R = 0$. It has a resonance at this frequency, and because there is no damping, the response at resonance is infinite. Of course, any real LRC circuit will have some damping, however small (cf. figure j on page 185).

Resonance with damping

example 36

- ▷ What is the amplitude of the current in a series LRC circuit?
- ▷ Generalizing from example 35, we add a third, real impedance:

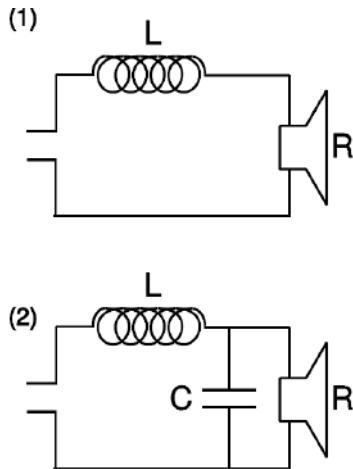
$$\begin{aligned} |\tilde{I}| &= \frac{|\tilde{V}|}{|Z|} \\ &= \frac{|\tilde{V}|}{|R + i\omega L - i/\omega C|} \\ &= \frac{|\tilde{V}|}{\sqrt{R^2 + (\omega L - 1/\omega C)^2}} \end{aligned}$$

This result would have taken pages of algebra without the complex number technique!

A second-order stereo crossover filter

example 37

A stereo crossover filter ensures that the high frequencies go to the tweeter and the lows to the woofer. This can be accomplished simply by putting a single capacitor in series with the tweeter and a single inductor in series with the woofer. However, such a filter does not cut off very sharply. Suppose we model the speakers as resistors. (They really have inductance as well, since they have coils in them that serve as electromagnets to move the diaphragm that makes the sound.) Then the power they draw is I^2R . Putting an inductor in series with the woofer, $\text{ag}/1$, gives a total impedance that at high frequencies is dominated by the inductor's, so the current is proportional to ω^{-1} , and the power drawn by the woofer is proportional to ω^{-2} .



ag / Example 37.

All the current passes through the inductor, so if the driving voltage being supplied on the left is \tilde{V}_d , we have

$$Z = Z_L + (Z_C^{-1} + Z_R^{-1})^{-1}$$

$$\tilde{V}_d = \tilde{I}_L Z,$$

and we also have

$$\tilde{V}_l = \tilde{l}_l Z_l$$

The loop rule, applied to the outer perimeter of the circuit, gives

$$\tilde{V}_d = \tilde{V}_L + \tilde{V}_R$$

Straightforward algebra now results in

$$\tilde{V}_R = \frac{\tilde{V}_d}{1 + Z_L/Z_C + Z_L/Z_B}.$$

At high frequencies, the Z_L/Z_C term, which varies as ω^2 , dominates, so \tilde{V}_R and \tilde{I}_R are proportional to ω^{-2} , and the power is proportional to ω^{-4} .

10.6 Fields by Gauss' law

10.6.1 Gauss' law

The flea of subsection 10.3.2 had a long and illustrious scientific career, and we're now going to pick up her story where we left off. This flea, whose name is Gauss⁹, has derived the equation $E_{\perp} = 2\pi k\sigma$ for the electric field very close to a charged surface with charge density σ . Next we will describe two improvements she is going to make to that equation.

First, she realizes that the equation is not as useful as it could be, because it only gives the part of the field *due to the surface*. If other charges are nearby, then their fields will add to this field as vectors, and the equation will not be true unless we carefully subtract out the field from the other charges. This is especially problematic for her because the planet on which she lives, known for obscure reasons as planet Flatcat, is itself electrically charged, and so are all the fleas — the only thing that keeps them from floating off into outer space is that they are negatively charged, while Flatcat carries a positive charge, so they are electrically attracted to it. When Gauss found the original version of her equation, she wanted to demonstrate it to her skeptical colleagues in the laboratory, using electric field meters and charged pieces of metal foil. Even if she set up the measurements by remote control, so that her the charge on her own body would be too far away to have any effect, they would be disrupted by the ambient field of planet Flatcat. Finally, however, she realized that she could improve her equation by rewriting it as follows:

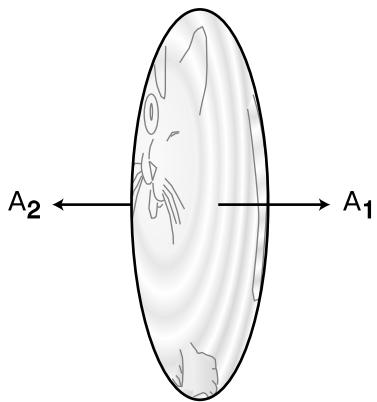
$$E_{\text{outward, on side 1}} + E_{\text{outward, on side 2}} = 4\pi k\sigma.$$

The tricky thing here is that “outward” means a different thing, depending on which side of the foil we’re on. On the left side, “outward” means to the left, while on the right side, “outward” is right. A positively charged piece of metal foil has a field that points leftward on the left side, and rightward on its right side, so the two contributions of $2\pi k\sigma$ are both positive, and we get $4\pi k\sigma$. On the other hand, suppose there is a field created by other charges, not by the charged foil, that happens to point to the right. On the right side, this externally created field is in the same direction as the foil’s field, but on the left side, the it *reduces* the strength of the leftward field created by the foil. The increase in one term of the equation balances the decrease in the other term. This new version of the equation is thus exactly correct regardless of what externally generated fields are present!

Her next innovation starts by multiplying the equation on both sides by the area, A , of one side of the foil:

$$(E_{\text{outward, on side 1}} + E_{\text{outward, on side 2}}) A = 4\pi k\sigma A$$

⁹no relation to the human mathematician of the same name



or

$$E_{\text{outward, on side } 1} A + E_{\text{outward, on side } 2} A = 4\pi kq,$$

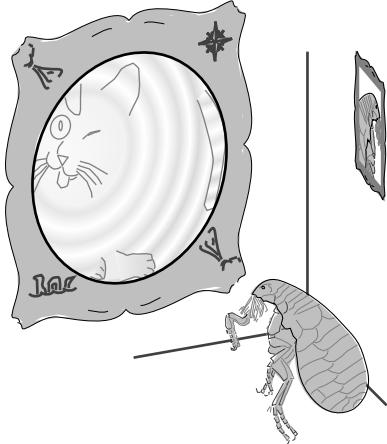
a / The area vector is defined to be perpendicular to the surface, in the outward direction. Its magnitude tells how much the area is.

where q is the charge of the foil. The reason for this modification is that she can now make the whole thing more attractive by defining a new vector, the area vector \mathbf{A} . As shown in figure a, she defines an area vector for side 1 which has magnitude A and points outward from side 1, and an area vector for side 2 which has the same magnitude and points outward from that side, which is in the opposite direction. The dot product of two vectors, $\mathbf{u} \cdot \mathbf{v}$, can be interpreted as $u_{\text{parallel to } v} |\mathbf{v}|$, and she can therefore rewrite her equation as

$$\mathbf{E}_1 \cdot \mathbf{A}_1 + \mathbf{E}_2 \cdot \mathbf{A}_2 = 4\pi kq.$$

The quantity on the left side of this equation is called the *flux* through the surface, written Φ .

Gauss now writes a grant proposal to her favorite funding agency, the BSGS (Blood-Suckers' Geological Survey), and it is quickly approved. Her audacious plan is to send out exploring teams to chart the electric fields of the whole planet of Flatcat, and thereby determine the total electric charge of the planet. The fleas' world is commonly assumed to be a flat disk, and its size is known to be finite, since the sun passes behind it at sunset and comes back around on the other side at dawn. The most daring part of the plan is that it requires surveying not just the known side of the planet but the uncharted Far Side as well. No flea has ever actually gone around the edge and returned to tell the tale, but Gauss assures them that they won't fall off — their negatively charged bodies will be attracted to the disk no matter which side they are on.



b / Gauss contemplates map of the known world.

a

Of course it is possible that the electric charge of planet Flatcat is not perfectly uniform, but that isn't a problem. As discussed in subsection 10.3.2, as long as one is very close to the surface, the field only depends on the *local* charge density. In fact, a side-benefit of Gauss's program of exploration is that any such local irregularities will be mapped out. But what the newspapers find exciting is the idea that once all the teams get back from their voyages and tabulate their data, the *total* charge of the planet will have been determined for the first time. Each surveying team is assigned to visit a certain list of republics, duchies, city-states, and so on. They are to record each territory's electric field vector, as well as its area. Because the electric field may be nonuniform, the final equation for determining the planet's electric charge will have many terms, not just one for each side of the planet:

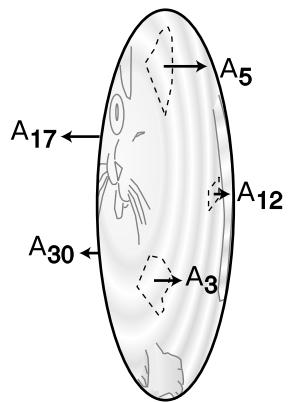
$$\Phi = \sum \mathbf{E}_j \cdot \mathbf{A}_j = 4\pi k q_{\text{total}}$$

Gauss herself leads one of the expeditions, which heads due east, toward the distant Tail Kingdom, known only from fables and the occasional account from a caravan of traders. A strange thing happens, however. Gauss embarks from her college town in the wetlands of the Tongue Republic, travels straight east, passes right through the Tail Kingdom, and one day finds herself right back at home, all without ever seeing the edge of the world! What can have happened? All at once she realizes that the world isn't flat.

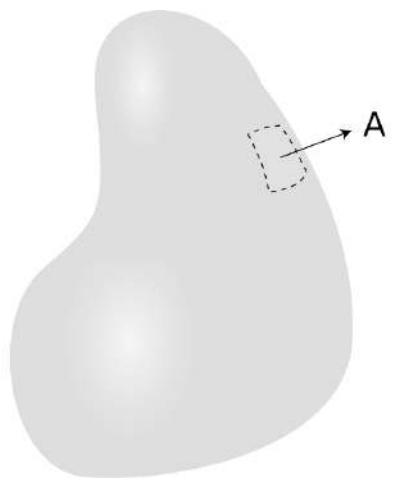
Now what? The surveying teams all return, the data are tabulated, and the result for the total charge of Flatcat is $(1/4\pi k) \sum \mathbf{E}_j \cdot \mathbf{A}_j = 37 \text{ nC}$ (units of nanocoulombs). But the equation was derived under the assumption that Flatcat was a disk. If Flatcat is really round, then the result may be completely wrong. Gauss and two of her grad students go to their favorite bar, and decide to keep on ordering Bloody Marys until they either solve their problems or forget them. One student suggests that perhaps Flatcat really is a disk, but the edges are rounded. Maybe the surveying teams really did flip over the edge at some point, but just didn't realize it. Under this assumption, the original equation will be approximately valid, and 37 nC really is the total charge of Flatcat.

A second student, named Newton, suggests that they take seriously the possibility that Flatcat is a sphere. In this scenario, their planet's surface is really curved, but the surveying teams just didn't notice the curvature, since they were close to the surface, and the surface was so big compared to them. They divided up the surface into a patchwork, and each patch was fairly small compared to the whole planet, so each patch was nearly flat. Since the patch is nearly flat, it makes sense to define an area vector that is perpendicular to it. In general, this is how we define the direction of an area vector, as shown in figure d. This only works if the areas are small. For instance, there would be no way to define an area vector for an entire sphere, since "outward" is in more than one direction.

If Flatcat is a sphere, then the inside of the sphere must be vast, and there is no way of knowing exactly how the charge is arranged below the surface. However, the survey teams all found that the electric field was approximately perpendicular to the surface everywhere, and that its strength didn't change very much from one location to another. The simplest explanation is that the charge is all concentrated in one small lump at the center of the sphere. They have no way of knowing if this is really the case, but it's a hypothesis that allows them to see how much their 37 nC result would change if they assumed a different geometry. Making this assumption, Newton performs the following simple computation on a napkin. The field at the surface is related to the charge at the center by



c / Each part of the surface has its own area vector. Note the differences in lengths of the vectors, corresponding to the unequal areas.



d / An area vector can be defined for a sufficiently small part of a curved surface.

$$|\mathbf{E}| = \frac{kq_{total}}{r^2},$$

where r is the radius of Flatcat. The flux is then

$$\Phi = \sum \mathbf{E}_j \cdot \mathbf{A}_j,$$

and since the \mathbf{E}_j and \mathbf{A}_j vectors are parallel, the dot product equals $|\mathbf{E}_j||\mathbf{A}_j|$, so

$$\Phi = \sum \frac{kq_{total}}{r^2} |\mathbf{A}_j|.$$

But the field strength is always the same, so we can take it outside the sum, giving

$$\begin{aligned}\Phi &= \frac{kq_{total}}{r^2} \sum |\mathbf{A}_j| \\ &= \frac{kq_{total}}{r^2} A_{total} \\ &= \frac{kq_{total}}{r^2} 4\pi r^2 \\ &= 4\pi kq_{total}.\end{aligned}$$

Not only have all the factors of r canceled out, but the result is the same as for a disk!

Everyone is pleasantly surprised by this apparent mathematical coincidence, but is it anything more than that? For instance, what if the charge wasn't concentrated at the center, but instead was evenly distributed throughout Flatcat's interior volume? Newton, however, is familiar with a result called the shell theorem (page 102), which states that the field of a uniformly charged sphere is the same as if all the charge had been concentrated at its center.¹⁰ We now have three different assumptions about the shape of Flatcat and the arrangement of the charges inside it, and all three lead to exactly the *same* mathematical result, $\Phi = 4\pi kq_{total}$. This is starting to look like more than a coincidence. In fact, there is a general mathematical theorem, called Gauss' theorem, which states the following:

For any region of space, the flux through the surface equals $4\pi kq_{in}$, where q_{in} is the total charge in that region.

Don't memorize the factor of 4π in front — you can rederive it any time you need to, by considering a spherical surface centered on a point charge.

¹⁰Newton's human namesake actually proved this for gravity, not electricity, but they're both $1/r^2$ forces, so the proof works equally well in both cases.

Note that although region and its surface had a definite physical existence in our story — they are the planet Flatcat and the surface of planet Flatcat — Gauss' law is true for any region and surface we choose, and in general, the Gaussian surface has no direct physical significance. It's simply a computational tool.

Rather than proving Gauss' theorem and then presenting some examples and applications, it turns out to be easier to show some examples that demonstrate its salient properties. Having understood these properties, the proof becomes quite simple.

self-check K

Suppose we have a negative point charge, whose field points inward, and we pick a Gaussian surface which is a sphere centered on that charge. How does Gauss' theorem apply here? \triangleright Answer, p. 1064

10.6.2 Additivity of flux

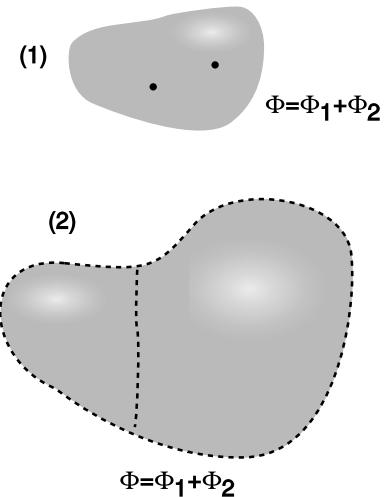
Figure e shows two two different ways in which flux is additive. Figure e/1, additivity by charge, shows that we can break down a charge distribution into two or more parts, and the flux equals the sum of the fluxes due to the individual charges. This follows directly from the fact that the flux is defined in terms of a dot product, $\mathbf{E} \cdot \mathbf{A}$, and the dot product has the additive property $(\mathbf{a} + \mathbf{b}) \cdot \mathbf{c} = \mathbf{a} \cdot \mathbf{c} + \mathbf{b} \cdot \mathbf{c}$.

To understand additivity of flux by region, e/2, we have to consider the parts of the two surfaces that were eliminated when they were joined together, like knocking out a wall to make two small apartments into one big one. Although the two regions shared this wall before it was removed, the area vectors were opposite: the direction that is outward from one region is inward with respect to the other. Thus if the field on the wall contributes positive flux to one region, it contributes an equal amount of negative flux to the other region, and we can therefore eliminate the wall to join the two regions, without changing the total flux.

10.6.3 Zero flux from outside charges

A third important property of Gauss' theorem is that it only refers to the charge *inside* the region we choose to discuss. In other words, it asserts that any charge outside the region contributes zero to the flux. This makes at least some sense, because a charge outside the region will have field vectors pointing into the surface on one side, and out of the surface on the other. Certainly there should be at least partial cancellation between the negative (inward) flux on one side and the positive (outward) flux on the other. But why should this cancellation be exact?

To see the reason for this perfect cancellation, we can imagine space as being built out of tiny cubes, and we can think of any charge distribution as being composed of point charges. The additivity-by-charge property tells us that any charge distribution can be handled



e / 1. The flux due to two charges equals the sum of the fluxes from each one. 2. When two regions are joined together, the flux through the new region equals the sum of the fluxes through the two parts.

by considering its point charges individually, and the additivity-by-region property tells us that if we have a single point charge outside a big region, we can break the region down into tiny cubes. If we can prove that the flux through such a tiny cube really does cancel exactly, then the same must be true for any region, which we could build out of such cubes, and any charge distribution, which we can build out of point charges.

For simplicity, we will carry out this calculation only in the special case shown in figure f, where the charge lies along one axis of the cube. Let the sides of the cube have length $2b$, so that the area of each side is $(2b)^2 = 4b^2$. The cube extends a distance b above, below, in front of, and behind the horizontal x axis. There is a distance $d - b$ from the charge to the left side, and $d + b$ to the right side.

There will be one negative flux, through the left side, and five positive ones. Of these positive ones, the one through the right side is very nearly the same in magnitude as the negative flux through the left side, but just a little less because the field is weaker on the right, due to the greater distance from the charge. The fluxes through the other four sides are very small, since the field is nearly perpendicular to their area vectors, and the dot product $\mathbf{E}_j \cdot \mathbf{A}_j$ is zero if the two vectors are perpendicular. In the limit where b is very small, we can approximate the flux by evaluating the field at the center of each of the cube's six sides, giving

$$\begin{aligned}\Phi &= \Phi_{left} + 4\Phi_{side} + \Phi_{right} \\ &= |\mathbf{E}_{left}| |\mathbf{A}_{left}| \cos 180^\circ + 4|\mathbf{E}_{side}| |\mathbf{A}_{side}| \cos \theta_{side} \\ &\quad + |\mathbf{E}_{right}| |\mathbf{A}_{right}| \cos 0^\circ,\end{aligned}$$

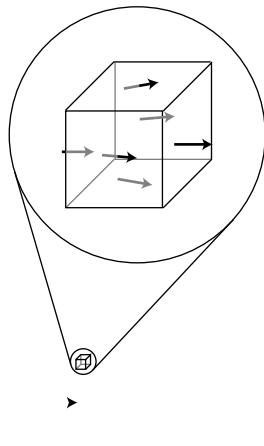
and a little trig gives $\cos \theta_{side} \approx b/d$, so

$$\begin{aligned}\Phi &= -|\mathbf{E}_{left}| |\mathbf{A}_{left}| + 4|\mathbf{E}_{side}| |\mathbf{A}_{side}| \frac{b}{d} + |\mathbf{E}_{right}| |\mathbf{A}_{right}| \\ &= (4b^2) \left(-|\mathbf{E}_{left}| + 4|\mathbf{E}_{side}| \frac{b}{d} + |\mathbf{E}_{right}| \right) \\ &= (4b^2) \left(-\frac{kq}{(d-b)^2} + 4 \frac{kq}{d^2} \frac{b}{d} + \frac{kq}{(d+b)^2} \right) \\ &= \left(\frac{4kqb^2}{d^2} \right) \left(-\frac{1}{(1-b/d)^2} + \frac{4b}{d} + \frac{1}{(1+b/d)^2} \right).\end{aligned}$$

Using the approximation $(1+\epsilon)^{-2} \approx 1-2\epsilon$ for small ϵ , this becomes

$$\begin{aligned}\Phi &= \left(\frac{4kqb^2}{d^2} \right) \left(-1 - \frac{2b}{d} + \frac{4b}{d} + 1 - \frac{2b}{d} \right) \\ &= 0.\end{aligned}$$

Thus in the limit of a very small cube, $b \ll d$, we have proved that the flux due to this exterior charge is zero. The proof can be



f / The flux through a tiny cube due to a point charge.

extended to the case where the charge is not along any axis of the cube,¹¹ and based on additivity we then have a proof that the flux due to an outside charge is always zero.

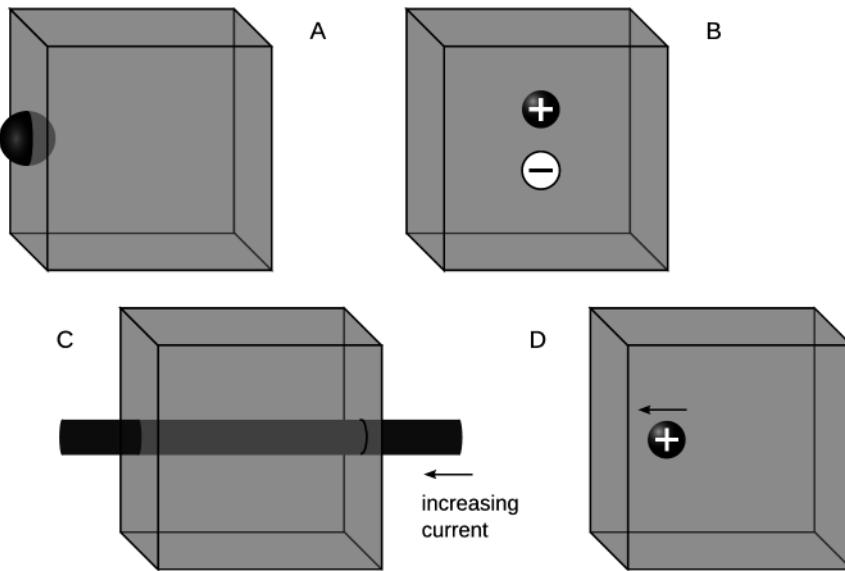
No charge on the interior of a conductor *example 38*

I asserted on p. 545 that for a perfect conductor in equilibrium, excess charge is found only at the surface, never in the interior. This can be proved using Gauss's theorem. Suppose that a charge q existed at some point in the interior, and it was in stable equilibrium. For concreteness, let's say q is positive. If its equilibrium is to be stable, then we need an electric field everywhere around it that points inward like a pincushion, so that if the charge were to be perturbed slightly, the field would bring it back to its equilibrium position. Since Newton's third law forbids objects from making forces on themselves, this field would have to be the field contributed by all the other charges, not by q itself. But this is impossible, because this kind of inward-pointing pincushion pattern would have a nonzero (negative) flux through the pincushion, but Gauss's theorem says we can't have flux from outside charges.

¹¹The math gets messy for the off-axis case. This part of the proof can be completed more easily and transparently using the techniques of section 10.7, and that is exactly we'll do in example 40 on page 655.

Discussion Questions

g / Discussion question A-D.



A One question that might naturally occur to you about Gauss's law is what happens for charge that is exactly on the surface — should it be counted toward the enclosed charge, or not? If charges can be perfect, infinitesimal points, then this could be a physically meaningful question. Suppose we approach this question by way of a limit: start with charge q spread out over a sphere of finite size, and then make the size of the sphere approach zero. The figure shows a uniformly charged sphere that's exactly half-way in and half-way out of the cubical Gaussian surface. What is the flux through the cube, compared to what it would be if the charge was entirely enclosed? (There are at least three ways to find this flux: by direct integration, by Gauss's law, or by the additivity of flux by region.)

B The dipole is completely enclosed in the cube. What does Gauss's law say about the flux through the cube? If you imagine the dipole's field pattern, can you verify that this makes sense?

C The wire passes in through one side of the cube and out through the other. If the current through the wire is increasing, then the wire will act like an inductor, and there will be a voltage difference between its ends. (The inductance will be relatively small, since the wire isn't coiled up, and the ΔV will therefore also be fairly small, but still not zero.) The ΔV implies the existence of electric fields, and yet Gauss's law says the flux must be zero, since there is no charge inside the cube. Why isn't Gauss's law violated?

D The charge has been loitering near the edge of the cube, but is then suddenly hit with a mallet, causing it to fly off toward the left side of the cube. We haven't yet discussed in detail how disturbances in the electric and magnetic fields ripple outward through space, but it turns out that they do so at the speed of light. (In fact, that's what light is: ripples in the electric and magnetic fields.) Because the charge is closer to the left side of the cube, the change in the electric field occurs there before

the information reaches the right side. This would seem certain to lead to a violation of Gauss's law. How can the ideas explored in discussion question C show the resolution to this paradox?

10.6.4 Proof of Gauss' theorem

With the computational machinery we've developed, it is now simple to prove Gauss' theorem. Based on additivity by charge, it suffices to prove the law for a point charge. We have already proved Gauss' law for a point charge in the case where the point charge is outside the region. If we can prove it for the inside case, then we're all done.

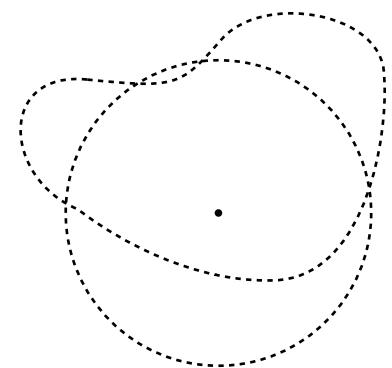
If the charge is inside, we reason as follows. First, we forget about the actual Gaussian surface of interest, and instead construct a spherical one, centered on the charge. For the case of a sphere, we've already seen the proof written on a napkin by the flea named Newton (page 643). Now wherever the actual surface sticks out beyond the sphere, we glue appropriately shaped pieces onto the sphere. In the example shown in figure h, we have to add two Mickey Mouse ears. Since these added pieces do not contain the point charge, the flux through them is zero, and additivity of flux by region therefore tells us that the total flux is not changed when we make this alteration. Likewise, we need to chisel out any regions where the sphere sticks out beyond the actual surface. Again, there is no change in flux, since the region being altered doesn't contain the point charge. This proves that the flux through the Gaussian surface of interest is the same as the flux through the sphere, and since we've already proved that that flux equals $4\pi kq_{in}$, our proof of Gauss' theorem is complete.

Discussion Questions

- A** A critical part of the proof of Gauss' theorem was the proof that a tiny cube has zero flux through it due to an external charge. Discuss qualitatively why this proof would fail if Coulomb's law was a $1/r$ or $1/r^3$ law.

10.6.5 Gauss' law as a fundamental law of physics

Note that the proof of Gauss' theorem depended on the computation on the napkin discussed on page 10.6.1. The crucial point in this computation was that the electric field of a point charge falls off like $1/r^2$, and since the area of a sphere is proportional to r^2 , the result is independent of r . The $1/r^2$ variation of the field also came into play on page 646 in the proof that the flux due to an outside charge is zero. In other words, if we discover some other force of nature which is proportional to $1/r^3$ or r , then Gauss' theorem will not apply to that force. Gauss' theorem is not true for nuclear forces, which fall off exponentially with distance. However, this is the *only* assumption we had to make about the nature of the field. Since gravity, for instance, also has fields that fall off as $1/r^2$, Gauss' theorem is equally valid for gravity — we just have to replace mass



h / Completing the proof of Gauss' theorem.

with charge, change the Coulomb constant k to the gravitational constant G , and insert a minus sign because the gravitational fields around a (positive) mass point inward.

Gauss' theorem can only be proved if we assume a $1/r^2$ field, and the converse is also true: any field that satisfies Gauss' theorem must be a $1/r^2$ field. Thus although we previously thought of Coulomb's law as the fundamental law of nature describing electric forces, it is equally valid to think of Gauss' theorem as the basic law of nature for electricity. From this point of view, Gauss' theorem is not a mathematical fact but an experimentally testable statement about nature, so we'll refer to it as *Gauss' law*, just as we speak of Coulomb's *law* or Newton's *law* of gravity.

If Gauss' law is equivalent to Coulomb's law, why not just use Coulomb's law? First, there are some cases where calculating a field is easy with Gauss' law, and hard with Coulomb's law. More importantly, Gauss' law and Coulomb's law are only mathematically equivalent under the assumption that all our charges are standing still, and all our fields are constant over time, i.e., in the study of electrostatics, as opposed to electrodynamics. As we broaden our scope to study generators, inductors, transformers, and radio antennas, we will encounter cases where Gauss' law is valid, but Coulomb's law is not.

10.6.6 Applications

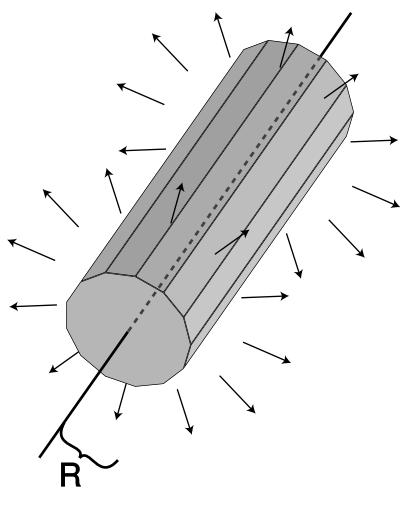
Often we encounter situations where we have a static charge distribution, and we wish to determine the field. Although superposition is a generic strategy for solving this type of problem, if the charge distribution is symmetric in some way, then Gauss' law is often a far easier way to carry out the computation.

Field of a long line of charge

Consider the field of an infinitely long line of charge, holding a uniform charge per unit length λ . Computing this field by brute-force superposition was fairly laborious (examples 13 on page 598 and 16 on page 605). With Gauss' law it becomes a very simple calculation.

The problem has two types of symmetry. The line of charge, and therefore the resulting field pattern, look the same if we rotate them about the line. The second symmetry occurs because the line is infinite: if we slide the line along its own length, nothing changes. This sliding symmetry, known as a translation symmetry, tells us that the field must point directly away from the line at any given point.

Based on these symmetries, we choose the Gaussian surface shown in figure i. If we want to know the field at a distance R from the line, then we choose this surface to have a radius R , as



i / Applying Gauss' law to an infinite line of charge.

shown in the figure. The length, L , of the surface is irrelevant.

The field is parallel to the surface on the end caps, and therefore perpendicular to the end caps' area vectors, so there is no contribution to the flux. On the long, thin strips that make up the rest of the surface, the field is perpendicular to the surface, and therefore parallel to the area vector of each strip, so that the dot product occurring in the definition of the flux is $\mathbf{E}_j \cdot \mathbf{A}_j = |\mathbf{E}_j||\mathbf{A}_j|\cos 0^\circ = |\mathbf{E}_j||\mathbf{A}_j|$. Gauss' law gives

$$4\pi k q_{in} = \sum \mathbf{E}_j \cdot \mathbf{A}_j$$

$$4\pi k \lambda L = \sum |\mathbf{E}_j||\mathbf{A}_j|.$$

The magnitude of the field is the same on every strip, so we can take it outside the sum.

$$4\pi k \lambda L = |\mathbf{E}| \sum |\mathbf{A}_j|$$

In the limit where the strips are infinitely narrow, the surface becomes a cylinder, with (area)=(circumference)(length)= $2\pi RL$.

$$4\pi k \lambda L = |\mathbf{E}| \times 2\pi RL$$

$$|\mathbf{E}| = \frac{2k\lambda}{R}$$

Field near a surface charge

As claimed earlier, the result $E = 2\pi k\sigma$ for the field near a charged surface is a special case of Gauss' law. We choose a Gaussian surface of the shape shown in figure j, known as a Gaussian pillbox. The exact shape of the flat end caps is unimportant.

The symmetry of the charge distribution tells us that the field points directly away from the surface, and is equally strong on both sides of the surface. This means that the end caps contribute equally to the flux, and the curved sides have zero flux through them. If the area of each end cap is A , then

$$4\pi k q_{in} = \mathbf{E}_1 \cdot \mathbf{A}_1 + \mathbf{E}_2 \cdot \mathbf{A}_2,$$

where the subscripts 1 and 2 refer to the two end caps. We have $\mathbf{A}_2 = -\mathbf{A}_1$, so

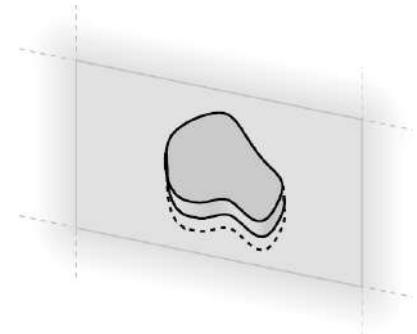
$$4\pi k q_{in} = \mathbf{E}_1 \cdot \mathbf{A}_1 - \mathbf{E}_2 \cdot \mathbf{A}_1$$

$$4\pi k q_{in} = (\mathbf{E}_1 - \mathbf{E}_2) \cdot \mathbf{A}_1,$$

and by symmetry the magnitudes of the two fields are equal, so

$$2|\mathbf{E}|A = 4\pi k\sigma A$$

$$|\mathbf{E}| = 2\pi k\sigma$$



j / Applying Gauss' law to an infinite charged surface.

The symmetry between the two sides could be broken by the existence of other charges nearby, whose fields would add onto the field of the surface itself. Even then, Gauss's law still guarantees

$$4\pi k q_{in} = (\mathbf{E}_1 - \mathbf{E}_2) \cdot \mathbf{A}_1,$$

or

$$|\mathbf{E}_{\perp,1} - \mathbf{E}_{\perp,2}| = 4\pi k \sigma,$$

where the subscript \perp indicates the component of the field parallel to the surface (i.e., parallel to the area vectors). In other words, the electric field changes discontinuously when we pass through a charged surface; the discontinuity occurs in the component of the field perpendicular to the surface, and the amount of discontinuous change is $4\pi k \sigma$. This is a completely general statement that is true near any charged surface, regardless of the existence of other charges nearby.

10.7 Gauss' law in differential form

10.7.1 Gauss's law as a local law

Gauss' law is a bit spooky. It relates the field on the Gaussian surface to the charges inside the surface. What if the charges have been moving around, and the field at the surface right now is the one that was created by the charges in their previous locations? Gauss' law — unlike Coulomb's law — still works in cases like these, but it's far from obvious how the flux and the charges can still stay in agreement if the charges have been moving around.

For this reason, it would be more physically attractive to restate Gauss' law in a different form, so that it related the behavior of the field at one point to the charges that were actually present at that point. This is essentially what we were doing in the fable of the flea named Gauss: the fleas' plan for surveying their planet was essentially one of dividing up the surface of their planet (which they believed was flat) into a patchwork, and then constructing *small* Gaussian pillbox around each *small* patch. The equation $E_{\perp} = 2\pi k\sigma$ then related a particular property of the *local* electric field to the *local* charge density.

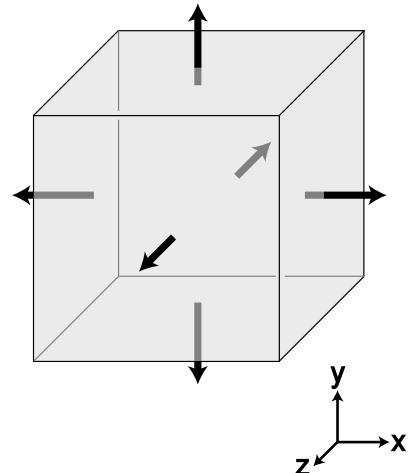
In general, charge distributions need not be confined to a flat surface — life is three-dimensional — but the general approach of defining very small Gaussian surfaces is still a good one. Our strategy is to divide up space into tiny cubes, like the one on page 645. Each such cube constitutes a Gaussian surface, which may contain some charge. Again we approximate the field using its six values at the center of each of the six sides. Let the cube extend from x to $x + dx$, from y to $y + dy$, and from z to $z + dz$.

The sides at x and $x + dx$ have area vectors $-dy dz \hat{x}$ and $dy dz \hat{x}$, respectively. The flux through the side at x is $-E_x(x) dy dz$, and the flux through the opposite side, at $x + dx$ is $E_x(x + dx) dy dz$. The sum of these is $(E_x(x + dx) - E_x(x)) dy dz$, and if the field was uniform, the flux through these two opposite sides would be zero. It will only be zero if the field's x component changes as a function of x . The difference $E_x(x + dx) - E_x(x)$ can be rewritten as $dE_x = (dE_x)/(dx) dx$, so the contribution to the flux from these two sides of the cube ends up being

$$\frac{dE_x}{dx} dx dy dz.$$

Doing the same for the other sides, we end up with a total flux

$$\begin{aligned} d\Phi &= \left(\frac{dE_x}{dx} + \frac{dE_y}{dy} + \frac{dE_z}{dz} \right) dx dy dz \\ &= \left(\frac{dE_x}{dx} + \frac{dE_y}{dy} + \frac{dE_z}{dz} \right) dv, \end{aligned}$$



a / A tiny cubical Gaussian surface.

where dv is the volume of the cube. In evaluating each of these three derivatives, we are going to treat the other two variables as constants, to emphasize this we use the partial derivative notation ∂ introduced in chapter 3,

$$d\Phi = \left(\frac{\partial E_x}{\partial x} + \frac{\partial E_y}{\partial y} + \frac{\partial E_z}{\partial z} \right) dv.$$

Using Gauss' law,

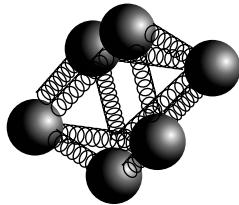
$$4\pi k q_{in} = \left(\frac{\partial E_x}{\partial x} + \frac{\partial E_y}{\partial y} + \frac{\partial E_z}{\partial z} \right) dv,$$

and we introduce the notation ρ (Greek letter rho) for the charge per unit volume, giving

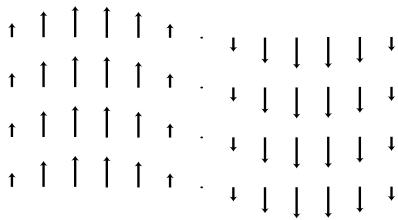
$$4\pi k \rho = \frac{\partial E_x}{\partial x} + \frac{\partial E_y}{\partial y} + \frac{\partial E_z}{\partial z}.$$

The quantity on the right is called the *divergence* of the electric field, written $\text{div } \mathbf{E}$. Using this notation, we have

$$\text{div } \mathbf{E} = 4\pi k \rho.$$



b / A meter for measuring $\text{div } \mathbf{E}$.



c / Example 39.

This equation has all the same physical implications as Gauss' law. After all, we proved Gauss' law by breaking down space into little cubes like this. We therefore refer to it as the differential form of Gauss' law, as opposed to $\Phi = 4\pi k q_{in}$, which is called the integral form.

Figure b shows an intuitive way of visualizing the meaning of the divergence. The meter consists of some electrically charged balls connected by springs. If the divergence is positive, then the whole cluster will expand, and it will contract its volume if it is placed at a point where the field has $\text{div } \mathbf{E} < 0$. What if the field is constant? We know based on the definition of the divergence that we should have $\text{div } \mathbf{E} = 0$ in this case, and the meter does give the right result: all the balls will feel a force in the same direction, but they will neither expand nor contract.

Divergence of a sine wave

example 39

▷ Figure c shows an electric field that varies as a sine wave. This is in fact what you'd see in a light wave: light is a wave pattern made of electric and magnetic fields. (The magnetic field would look similar, but would be in a plane perpendicular to the page.) What is the divergence of such a field, and what is the physical significance of the result?

▷ Intuitively, we can see that no matter where we put the div-meter in this field, it will neither expand nor contract. For instance, if we put it at the center of the figure, it will start spinning, but that's it.

Mathematically, let the x axis be to the right and let y be up. The field is of the form

$$\mathbf{E} = (\sin Kx) \hat{\mathbf{y}},$$

where the constant K is not to be confused with Coulomb's constant. Since the field has only a y component, the only term in the divergence we need to evaluate is

$$\mathbf{E} = \frac{\partial E_y}{\partial y},$$

but this vanishes, because E_y depends only on x , not y : we treat y as a constant when evaluating the partial derivative $\partial E_y / \partial y$, and the derivative of an expression containing only constants must be zero.

Physically this is a very important result: it tells us that a light wave can exist without any charges along the way to "keep it going." In other words, light can travel through a vacuum, a region with no particles in it. If this wasn't true, we'd be dead, because the sun's light wouldn't be able to get to us through millions of kilometers of empty space!

Electric field of a point charge

example 40

The case of a point charge is tricky, because the field behaves badly right on top of the charge, blowing up and becoming discontinuous. At this point, we cannot use the component form of the divergence, since none of the derivatives are well defined. However, a little visualization using the original definition of the divergence will quickly convince us that $\operatorname{div} \mathbf{E}$ is infinite here, and that makes sense, because the density of charge has to be infinite at a point where there is a zero-size point of charge (finite charge in zero volume).

At all other points, we have

$$\mathbf{E} = \frac{kq}{r^2} \hat{\mathbf{r}},$$

where $\hat{\mathbf{r}} = \mathbf{r}/r = (x\hat{\mathbf{x}} + y\hat{\mathbf{y}} + z\hat{\mathbf{z}})/r$ is the unit vector pointing radially away from the charge. The field can therefore be written as

$$\begin{aligned} \mathbf{E} &= \frac{kq}{r^3} \hat{\mathbf{r}} \\ &= \frac{kq(x\hat{\mathbf{x}} + y\hat{\mathbf{y}} + z\hat{\mathbf{z}})}{(x^2 + y^2 + z^2)^{3/2}}. \end{aligned}$$

The three terms in the divergence are all similar, e.g.,

$$\begin{aligned} \frac{\partial E_x}{\partial x} &= kq \frac{\partial}{\partial x} \left[\frac{x}{(x^2 + y^2 + z^2)^{3/2}} \right] \\ &= kq \left[\frac{1}{(x^2 + y^2 + z^2)^{3/2}} - \frac{3}{2} \frac{2x^2}{(x^2 + y^2 + z^2)^{5/2}} \right] \\ &= kq \left(r^{-3} - 3x^2 r^{-5} \right). \end{aligned}$$

Straightforward algebra shows that adding in the other two terms results in zero, which makes sense, because there is no charge except at the origin.

Gauss' law in differential form lends itself most easily to finding the charge density when we are given the field. What if we want to find the field given the charge density? As demonstrated in the following example, one technique that often works is to guess the general form of the field based on experience or physical intuition, and then try to use Gauss' law to find what specific version of that general form will be a solution.

The field inside a uniform sphere of charge *example 41*

- ▷ Find the field inside a uniform sphere of charge whose charge density is ρ . (This is very much like finding the gravitational field at some depth below the surface of the earth.)
- ▷ By symmetry we know that the field must be purely radial (in and out). We guess that the solution might be of the form

$$\mathbf{E} = br^\rho \hat{\mathbf{r}},$$

where r is the distance from the center, and b and ρ are constants. A negative value of ρ would indicate a field that was strongest at the center, while a positive ρ would give zero field at the center and stronger fields farther out. Physically, we know by symmetry that the field is zero at the center, so we expect ρ to be positive.

As in the example 40, we rewrite $\hat{\mathbf{r}}$ as \mathbf{r}/r , and to simplify the writing we define $n = \rho - 1$, so

$$\mathbf{E} = br^n \mathbf{r}.$$

Gauss' law in differential form is

$$\operatorname{div} \mathbf{E} = 4\pi k\rho,$$

so we want a field whose divergence is constant. For a field of the form we guessed, the divergence has terms in it like

$$\begin{aligned}\frac{\partial E_x}{\partial x} &= \frac{\partial}{\partial x} (br^n x) \\ &= b \left(nr^{n-1} \frac{\partial r}{\partial x} x + r^n \right)\end{aligned}$$

The partial derivative $\partial r/\partial x$ is easily calculated to be x/r , so

$$\frac{\partial E_x}{\partial x} = b \left(nr^{n-2} x^2 + r^n \right)$$

Adding in similar expressions for the other two terms in the divergence, and making use of $x^2 + y^2 + z^2 = r^2$, we have

$$\operatorname{div} \mathbf{E} = b(n+3)r^n.$$

This can indeed be constant, but only if n is 0 or -3 , i.e., ρ is 1 or -2 . The second solution gives a divergence which is constant and zero: this is the solution for the *outside* of the sphere! The first solution, which has the field directly proportional to r , must be the one that applies to the inside of the sphere, which is what we care about right now. Equating the coefficient in front to the one in Gauss' law, the field is

$$\mathbf{E} = \frac{4\pi k\rho}{3} r \hat{\mathbf{r}}.$$

The field is zero at the center, and gets stronger and stronger as we approach the surface.

Discussion Questions

- A** As suggested by the figure, discuss the results you would get by inserting the div-meter at various locations in the sine-wave field.

10.7.2 Poisson's equation and Laplace's equation

Gauss's law, $\text{div } \mathbf{E} = 4\pi k\rho$, can also be stated in terms of the potential. Since $\mathbf{E} = \nabla V$, we have $\text{div } \nabla V = 4\pi k\rho$. If we work out the combination of operators $\text{div } \nabla$ in a Cartesian coordinate system, we get $\partial^2/\partial x^2 + \partial^2/\partial y^2 + \partial^2/\partial z^2$, which is called the Laplacian and notated ∇^2 . The Laplacian is discussed in more detail on p. 912. The version of Gauss's law written in terms of the potential,

$$\nabla^2 V = 4\pi k\rho,$$

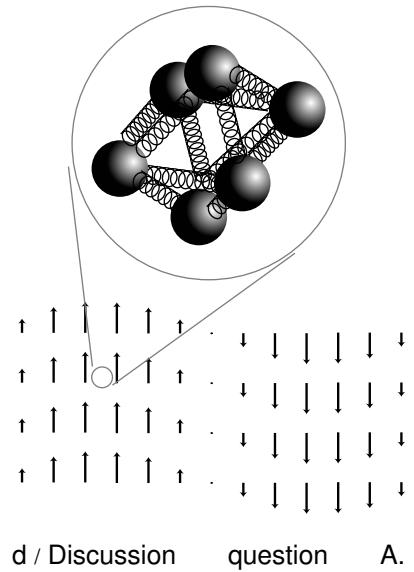
is called Poisson's equation, while in the special case of a vacuum, with $\rho = 0$, we have

$$\nabla^2 V = 0,$$

known as Laplace's equation. Many problems in electrostatics can be stated in terms of finding potential that satisfies Laplace's equation, usually with some set of *boundary conditions*. For example, if an infinite parallel-plate capacitor has plates parallel to the x - y plane at certain given potentials, then these plates form a boundary for the region between the plates, and Laplace's equation has a solution in this region of the form $V = az + b$. It's easy to verify that this is a solution of Laplace's equation, since all three of the partial derivatives vanish.

10.7.3 The method of images

A car's radio antenna is usually in the form of a whip sticking up above its metal roof. This is an example involving radio waves, which are time-varying electric and magnetic fields, but a similar, simpler electrostatic example is the following. Suppose that we position a charge $q > 0$ at a distance ℓ from a conducting plane. What is the resulting electric field? The conductor has charges that are free to move, and due to the field of the charge q , we will end up with a net concentration of negative charge in the part of the plane near

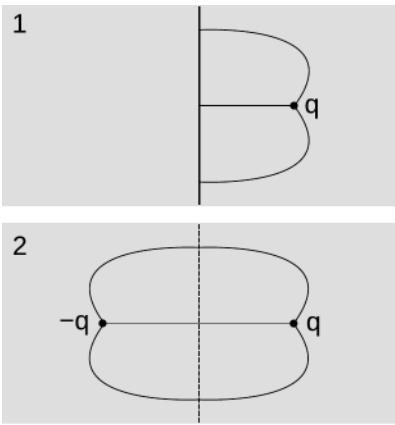


d / Discussion question A.

q . The field in the vacuum surrounding q will be a sum of fields due to q and fields due to these charges in the conducting plane. The problem can be stated as that of finding a solution to Poisson's equation with the boundary condition that $V = 0$ at the conducting plane. Figure e/1 shows the kind of field lines we expect.

This looks like a very complicated problem, but there is trick that allows us to find a simple solution. We can convert the problem into an equivalent one in which the conductor isn't present, but a fictitious *image* charge $-q$ is placed at an equal distance behind the plane, like a reflection in a mirror, as in figure e/2. The field is then simply the sum of the fields of the charges q and $-q$, so we can either add the field vectors or add the potentials. By symmetry, the field lines are perpendicular to the plane, so the plane is an surface of constant potential, as required.

This chapter is summarized on page 1086. Notation and terminology are tabulated on pages 1070-1071.



e / The method of images.

Problems

The symbols \checkmark , \blacksquare , etc. are explained on page 670.

1 The gap between the electrodes in an automobile engine's spark plug is 0.060 cm. To produce an electric spark in a gasoline-air mixture, an electric field of 3.0×10^6 V/m must be achieved. On starting a car, what minimum voltage must be supplied by the ignition circuit? Assume the field is uniform. \checkmark

(b) The small size of the gap between the electrodes is inconvenient because it can get blocked easily, and special tools are needed to measure it. Why don't they design spark plugs with a wider gap?



2 (a) As suggested in example 12 on page 597, use approximations to show that the expression given for the electric field approaches kQ/d^2 for large d .

(b) Do the same for the result of example 15 on page 601. \blacksquare

3 Astronomers believe that the mass distribution (mass per unit volume) of some galaxies may be approximated, in spherical coordinates, by $\rho = ae^{-br}$, for $0 \leq r \leq \infty$, where ρ is the density. Find the total mass.



4 (a) At time $t = 0$, a positively charged particle is placed, at rest, in a vacuum, in which there is a uniform electric field of magnitude E . Write an equation giving the particle's speed, v , in terms of t , E , and its mass and charge m and q . \checkmark

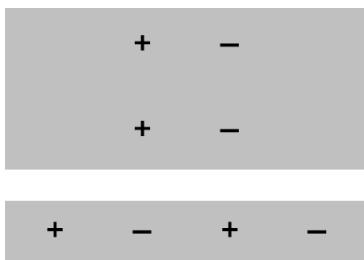
(b) If this is done with two different objects and they are observed to have the same motion, what can you conclude about their masses and charges? (For instance, when radioactivity was discovered, it was found that one form of it had the same motion as an electron in this type of experiment.) \blacksquare

5 Show that the alternative definition of the magnitude of the electric field, $|\mathbf{E}| = \tau/D_t \sin \theta$, has units that make sense. \blacksquare

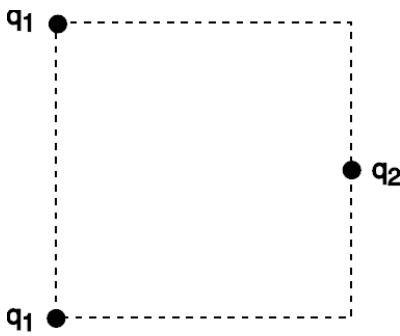
6 Redo the calculation of example 5 on page 588 using a different origin for the coordinate system, and show that you get the same result. \blacksquare

7 The definition of the dipole moment, $\mathbf{D} = \sum q_i \mathbf{r}_i$, involves the vector \mathbf{r}_i stretching from the origin of our coordinate system out to the charge q_i . There are clearly cases where this causes the dipole moment to be dependent on the choice of coordinate system. For instance, if there is only one charge, then we could make the dipole moment equal zero if we chose the origin to be right on top of the charge, or nonzero if we put the origin somewhere else.

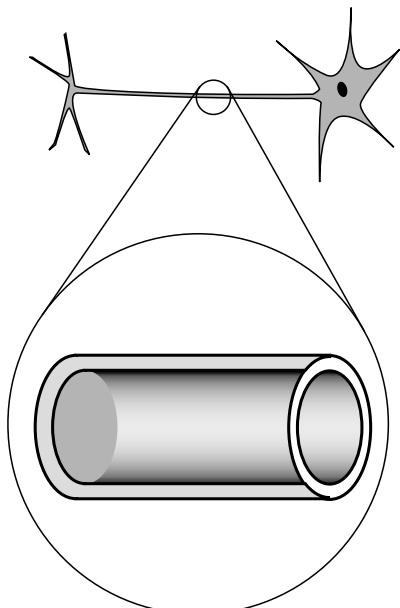
(a) Make up a numerical example with two charges of equal magnitude and opposite sign. Compute the dipole moment using two different coordinate systems that are oriented the same way, but



Problem 8.



Problem 11.



Problem 13.

differ in the choice of origin. Comment on the result.

(b) Generalize the result of part a to any pair of charges with equal magnitude and opposite sign. This is supposed to be a proof for *any* arrangement of the two charges, so don't assume any numbers.

(c) Generalize further, to n charges.

8 Compare the two dipole moments.

9 Find an arrangement of charges that has zero total charge and zero dipole moment, but that will make nonvanishing electric fields.

10 As suggested in example 14 on page 600, show that you can get the same result for the on-axis field by differentiating the potential.

11 Three charges are arranged on a square as shown. All three charges are positive. What value of q_2/q_1 will produce zero electric field at the center of the square? ✓

12 This is a one-dimensional problem, with everything confined to the x axis. Dipole A consists of a -1.000 C charge at $x = 0.000 \text{ m}$ and a 1.000 C charge at $x = 1.000 \text{ m}$. Dipole B has a -2.000 C charge at $x = 0.000 \text{ m}$ and a 2.000 C charge at $x = 0.500 \text{ m}$.

(a) Compare the two dipole moments.

(b) Calculate the field created by dipole A at $x = 10.000 \text{ m}$, and compare with the field dipole B would make. Comment on the result. ✓

13 In our by-now-familiar neuron, the voltage difference between the inner and outer surfaces of the cell membrane is about $V_{out} - V_{in} = -70 \text{ mV}$ in the resting state, and the thickness of the membrane is about 6.0 nm (i.e., only about a hundred atoms thick). What is the electric field inside the membrane? ✓

14 A proton is in a region in which the electric field is given by $E = a + bx^3$. If the proton starts at rest at $x_1 = 0$, find its speed, v , when it reaches position x_2 . Give your answer in terms of a , b , x_2 , and e and m , the charge and mass of the proton. ✓

15 (a) Given that the on-axis field of a dipole at large distances is proportional to D/r^3 , show that its potential varies as D/r^2 . (Ignore positive and negative signs and numerical constants of proportionality.)

(b) Write down an exact expression for the potential of a two-charge dipole at an on-axis point, without assuming that the distance is large compared to the size of the dipole. Your expression will have to contain the actual charges and size of the dipole, not just its dipole moment. Now use approximations to show that, at large distances, this is consistent with your answer to part a. ▷ Hint, p. 1036

16 A hydrogen atom is electrically neutral, so at large distances, we expect that it will create essentially zero electric field. This is

not true, however, near the atom or inside it. Very close to the proton, for example, the field is very strong. To see this, think of the electron as a spherically symmetric cloud that surrounds the proton, getting thinner and thinner as we get farther away from the proton. (Quantum mechanics tells us that this is a more correct picture than trying to imagine the electron orbiting the proton.) Near the center of the atom, the electron cloud's field cancels out by symmetry, but the proton's field is strong, so the total field is very strong. The potential in and around the hydrogen atom can be approximated using an expression of the form $V = r^{-1}e^{-r}$. (The units come out wrong, because I've left out some constants.) Find the electric field corresponding to this potential, and comment on its behavior at very large and very small r . \triangleright Solution, p. 1048 ■

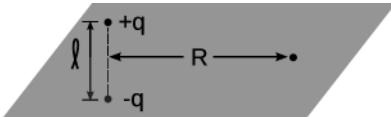
17 A carbon dioxide molecule is structured like O-C-O, with all three atoms along a line. The oxygen atoms grab a little bit of extra negative charge, leaving the carbon positive. The molecule's symmetry, however, means that it has no overall dipole moment, unlike a V-shaped water molecule, for instance. Whereas the potential of a dipole of magnitude D is proportional to D/r^2 , (see problem 15), it turns out that the potential of a carbon dioxide molecule at a distant point along the molecule's axis equals b/r^3 , where r is the distance from the molecule and b is a constant (cf. problem 9). What would be the electric field of a carbon dioxide molecule at a point on the molecule's axis, at a distance r from the molecule?

✓ ■

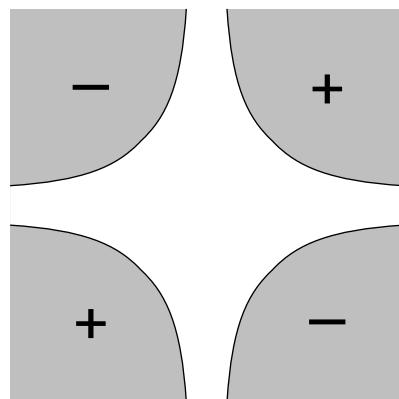
18 A hydrogen atom in a particular state has the charge density (charge per unit volume) of the electron cloud given by $\rho = ae^{-br}z^2$, where r is the distance from the proton, and z is the coordinate measured along the z axis. Given that the total charge of the electron cloud must be $-e$, find a in terms of the other variables. \checkmark ■

19 A dipole has a midplane, i.e., the plane that cuts through the dipole's center, and is perpendicular to the dipole's axis. Consider a two-charge dipole made of point charges $\pm q$ located at $z = \pm\ell/2$. Use approximations to find the field at a distant point in the midplane, and show that its magnitude comes out to be kD/R^3 (half what it would be at a point on the axis lying an equal distance from the dipole). ■

20 The figure shows a vacuum chamber surrounded by four metal electrodes shaped like hyperbolas. (Yes, physicists do sometimes ask their university machine shops for things machined in mathematical shapes like this. They have to be made on computer-controlled mills.) We assume that the electrodes extend far into and out of the page along the unseen z axis, so that by symmetry, the electric fields are the same for all z . The problem is therefore effectively two-dimensional. Two of the electrodes are at voltage $+V_0$, and the other two at $-V_0$, as shown. The equations of the hyperbolic surfaces are



Problem 19.



Problem 20.

$|xy| = b^2$, where b is a constant. (We can interpret b as giving the locations $x = \pm b$, $y = \pm b$ of the four points on the surfaces that are closest to the central axis.) There is no obvious, pedestrian way to determine the field or potential in the central vacuum region, but there's a trick that works: with a little mathematical insight, we see that the potential $V = V_0 b^{-2} xy$ is consistent with all the given information. (Mathematicians could prove that this solution was unique, but a physicist knows it on physical grounds: if there were two different solutions, there would be no physical way for the system to decide which one to do!)

(a) Use the techniques of subsection 10.2.2 to find the field in the vacuum region.

(b) Sketch the field as a “sea of arrows.”

✓ ■

21 (a) A certain region of three-dimensional space has a potential that varies as $V = br^2$, where r is the distance from the origin. Use the techniques of subsection 10.2.2 to find the field.

✓

(b) Write down another potential that gives exactly the same field.

■

22 (a) Example 13 on page 598 gives the field of a charged rod in its midplane. Starting from this result, take the limit as the length of the rod approaches infinity. Note that λ is not changing, so as L gets bigger, the total charge Q increases. ▷ Answer, p. 1069

(b) In the text, I have shown (by several different methods) that the field of an infinite, uniformly charged plane is $2\pi k\sigma$. Now you're going to rederive the same result by a different method. Suppose that it is the $x - y$ plane that is charged, and we want to find the field at the point $(0, 0, z)$. (Since the plane is infinite, there is no loss of generality in assuming $x = 0$ and $y = 0$.) Imagine that we slice the plane into an infinite number of straight strips parallel to the y axis. Each strip has infinitesimal width dx , and extends from x to $x + dx$. The contribution any one of these strips to the field at our point has a magnitude which can be found from part a. By vector addition, prove the desired result for the field of the plane of charge.

■

23 Consider the electric field created by a uniformly charged cylindrical surface that extends to infinity in one direction.

(a) Show that the field at the center of the cylinder's mouth is $2\pi k\sigma$, which happens to be the same as the field of an infinite *flat* sheet of charge!

(b) This expression is independent of the radius of the cylinder. Explain why this should be so. For example, what would happen if you doubled the cylinder's radius?

■

24 In an electrical storm, the cloud and the ground act like a parallel-plate capacitor, which typically charges up due to frictional electricity in collisions of ice particles in the cold upper atmosphere. Lightning occurs when the magnitude of the electric field reaches a



Problem 23.

critical value E_c , at which air is ionized.

(a) Treat the cloud as a flat square with sides of length L . If it is at a height h above the ground, find the amount of energy released in the lightning strike. ✓

(b) Based on your answer from part a, which is more dangerous, a lightning strike from a high-altitude cloud or a low-altitude one?

(c) Make an order-of-magnitude estimate of the energy released by a typical lightning bolt, assuming reasonable values for its size and altitude. E_c is about 10^6 N/C.

See problem 60 for a note on how recent research affects this estimate. ■

25 (a) Show that the energy in the electric field of a point charge is infinite! Does the integral diverge at small distances, at large distances, or both? ▷ Hint, p. 1036

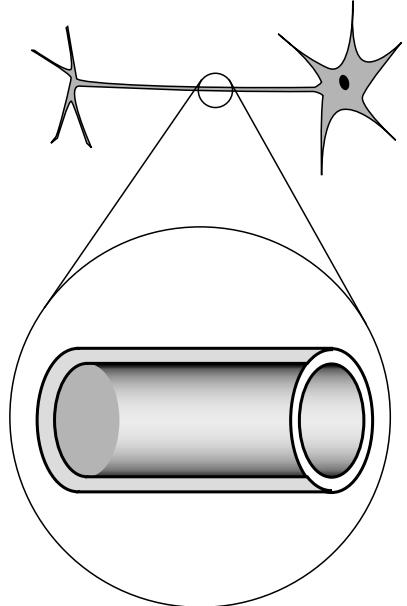
(b) Now calculate the energy in the electric field of a uniformly charged sphere with radius b . Based on the shell theorem, it can be shown that the field for $r > b$ is the same as for a point charge, while the field for $r < b$ is kqr/b^3 . (Example 41 shows this using a different technique.)

Remark: The calculation in part a seems to show that infinite energy would be required in order to create a charged, pointlike particle. However, there are processes that, for example, create electron-positron pairs, and these processes don't require infinite energy. According to Einstein's famous equation $E = mc^2$, the energy required to create such a pair should only be $2mc^2$, which is finite. ✓ ■

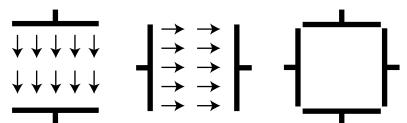
26 The neuron in the figure has been drawn fairly short, but some neurons in your spinal cord have tails (axons) up to a meter long. The inner and outer surfaces of the membrane act as the "plates" of a capacitor. (The fact that it has been rolled up into a cylinder has very little effect.) In order to function, the neuron must create a voltage difference V between the inner and outer surfaces of the membrane. Let the membrane's thickness, radius, and length be t , r , and L . (a) Calculate the energy that must be stored in the electric field for the neuron to do its job. (In real life, the membrane is made out of a substance called a dielectric, whose electrical properties increase the amount of energy that must be stored. For the sake of this analysis, ignore this fact.) ▷ Hint, p. 1036 ✓

(b) An organism's evolutionary fitness should be better if it needs less energy to operate its nervous system. Based on your answer to part a, what would you expect evolution to do to the dimensions t and r ? What other constraints would keep these evolutionary trends from going too far? ■

27 The figure shows cross-sectional views of two cubical capacitors, and a cross-sectional view of the same two capacitors put together so that their interiors coincide. A capacitor with the plates close together has a nearly uniform electric field between the plates,



Problem 26.



Problem 27.
Problems

and almost zero field outside; these capacitors don't have their plates very close together compared to the dimensions of the plates, but for the purposes of this problem, assume that they still have approximately the kind of idealized field pattern shown in the figure. Each capacitor has an interior volume of 1.00 m^3 , and is charged up to the point where its internal field is 1.00 V/m .

(a) Calculate the energy stored in the electric field of each capacitor when they are separate. ✓

(b) Calculate the magnitude of the interior field when the two capacitors are put together in the manner shown. Ignore effects arising from the redistribution of each capacitor's charge under the influence of the other capacitor. ✓

(c) Calculate the energy of the put-together configuration. Does assembling them like this release energy, consume energy, or neither? ✓ ■

28 Find the capacitance of the surface of the earth, assuming there is an outer spherical "plate" at infinity. (In reality, this outer plate would just represent some distant part of the universe to which we carried away some of the earth's charge in order to charge up the earth.) ✓ ■



Problem 29.

29 (a) Show that the field found in example 13 on page 598 reduces to $E = 2k\lambda/R$ in the limit of $L \rightarrow \infty$.

(b) An infinite strip of width b has a surface charge density σ . Find the field at a point at a distance z from the strip, lying in the plane perpendicularly bisecting the strip. ✓

(c) Show that this expression has the correct behavior in the limit where z approaches zero, and also in the limit of $z \gg b$. For the latter, you'll need the result of problem 22a, which is given on page 1069. ■

30 A solid cylinder of radius b and length ℓ is uniformly charged with a total charge Q . Find the electric field at a point at the center of one of the flat ends. ■

31 Find the potential at the edge of a uniformly charged disk. (Define $V = 0$ to be infinitely far from the disk.)

✓ ▷ Hint, p. 1036 ■

32 Find the energy stored in a capacitor in terms of its capacitance and the voltage difference across it. ✓ ■

33 (a) Find the capacitance of two identical capacitors in series.

(b) Based on this, how would you expect the capacitance of a parallel-plate capacitor to depend on the distance between the plates? ■

34 (a) Use complex number techniques to rewrite the function $f(t) = 4 \sin \omega t + 3 \cos \omega t$ in the form $A \sin(\omega t + \delta)$. ✓

(b) Verify the result using the trigonometric identity $\sin(\alpha + \beta) = \sin \alpha \cos \beta + \sin \beta \cos \alpha$. ■

35 (a) Show that the equation $V_L = L dI/dt$ has the right units.

- (b) Verify that RC has units of time.
(c) Verify that L/R has units of time. ■

36 Find the inductance of two identical inductors in parallel. ■

37 Calculate the quantity i^i (i.e., find its real and imaginary parts). ▷ Hint, p. 1037 ✓ ■

38 The wires themselves in a circuit can have resistance, inductance, and capacitance. Would “stray” inductance and capacitance be most important for low-frequency or for high-frequency circuits? For simplicity, assume that the wires act like they’re in *series* with an inductor or capacitor. ■

39 Starting from the relation $V = L dI/dt$ for the voltage difference across an inductor, show that an inductor has an impedance equal to $L\omega$. ■

40 A rectangular box is uniformly charged with a charge density ρ . The box is extremely long and skinny, and its cross-section is a square with sides of length b . The length is so great in comparison to b that we can consider it as being infinite. Find the electric field at a point lying on the box’s surface, at the midpoint between the two edges. Your answer will involve an integral that is most easily done using computer software. ■

41 A hollow cylindrical pipe has length ℓ and radius b . Its ends are open, but on the curved surface it has a charge density σ . A charge q with mass m is released at the center of the pipe, in unstable equilibrium. Because the equilibrium is unstable, the particle accelerates off in one direction or the other, along the axis of the pipe, and comes shooting out like a bullet from the barrel of a gun. Find the velocity of the particle when it’s infinitely far from the “gun.” Your answer will involve an integral that is difficult to do by hand; you may want to look it up in a table of integrals, do it online at integrals.com, or download and install the free Maxima symbolic math software from maxima.sourceforge.net. ■

42 Suppose that an FM radio tuner for the US commercial broadcast band (88-108 MHz) consists of a series LRC circuit. If the inductor is $1.0 \mu\text{H}$, what range of capacitances should the variable capacitor be able to provide? ✓ ■

43 (a) Find the parallel impedance of a $37 \text{ k}\Omega$ resistor and a 1.0 nF capacitor at $f = 1.0 \times 10^4 \text{ Hz}$. ✓

(b) A voltage with an amplitude of 1.0 mV drives this impedance at this frequency. What is the amplitude of the current drawn from the voltage source, what is the current’s phase angle with respect to the voltage, and does it lead the voltage, or lag behind it? ✓ ■

44 A series LRC circuit consists of a $1.000\ \Omega$ resistor, a $1.000\ F$ capacitor, and a $1.000\ H$ inductor. (These are not particularly easy values to find on the shelf!)

(a) Plot its impedance as a point in the complex plane for each of the following frequencies: $\omega=0.250, 0.500, 1.000, 2.000$, and $4.000\ Hz$.

(b) What is the resonant angular frequency, ω_{res} , and how does this relate to your plot? ✓

(c) What is the resonant frequency f_{res} corresponding to your answer in part b? ✓ ■

45 At a frequency of $53.1\ kHz$, a certain series LR circuit has an impedance of $1.6\ k\Omega + (1.2\ k\Omega)i$. Suppose that instead we want to achieve the same impedance using two circuit elements in parallel. What must the elements be? As a check on your answer, you should find that both values are round numbers when rounded off to the correct number of significant figures. ■

46 (a) Use Gauss' law to find the fields inside and outside an infinite cylindrical surface with radius b and uniform surface charge density σ . ✓

(b) Show that there is a discontinuity in the electric field equal to $4\pi k\sigma$ between one side of the surface and the other, as there should be (see page 652).

(c) Reexpress your result in terms of the charge per unit length, and compare with the field of a line of charge.

(d) A coaxial cable has two conductors: a central conductor of radius a , and an outer conductor of radius b . These two conductors are separated by an insulator. Although such a cable is normally used for time-varying signals, assume throughout this problem that there is simply a DC voltage between the two conductors. The outer conductor is thin, as in part c. The inner conductor is solid, but, as is always the case with a conductor in electrostatics, the charge is concentrated on the surface. Thus, you can find all the fields in part b by superposing the fields due to each conductor, as found in part c. (Note that on a given length of the cable, the total charge of the inner and outer conductors is zero, so $\lambda_1 = -\lambda_2$, but $\sigma_1 \neq \sigma_2$, since the areas are unequal.) Find the capacitance per unit length of such a cable. ✓ ■

47 In a certain region of space, the electric field is constant (i.e., the vector always has the same magnitude and direction). For simplicity, assume that the field points in the positive x direction.

(a) Use Gauss's law to prove that there is no charge in this region of space. This is most easily done by considering a Gaussian surface consisting of a rectangular box, whose edges are parallel to the x , y , and z axes.

(b) If there are no charges in this region of space, what could be making this electric field? ■

48 (a) In a series LC circuit driven by a DC voltage ($\omega = 0$), compare the energy stored in the inductor to the energy stored in the capacitor.

(b) Carry out the same comparison for an LC circuit that is oscillating freely (without any driving voltage).

(c) Now consider the general case of a series LC circuit driven by an oscillating voltage at an arbitrary frequency. Let $\overline{U_L}$ and $\overline{U_C}$ be the average energy stored in the inductor, and similarly for $\overline{U_C}$. Define a quantity $u = \overline{U_C}/(\overline{U_L} + \overline{U_C})$, which can be interpreted as the capacitor's average share of the energy, while $1 - u$ is the inductor's average share. Find u in terms of L , C , and ω , and sketch a graph of u and $1 - u$ versus ω . What happens at resonance? Make sure your result is consistent with your answer to part a. ✓



49 (a) Use Gauss' law to find the field inside an infinite cylinder with radius b and uniform charge density ρ . (The external field has the same form as the one in problem 46.) ✓

(b) Check that your answer makes sense on the axis.

(c) Check that the units of your answer make sense.



50 (a) In a certain region of space, the electric field is given by $\mathbf{E} = bx\hat{\mathbf{x}}$, where b is a constant. Find the amount of charge contained within a cubical volume extending from $x = 0$ to $x = a$, from $y = 0$ to $y = a$, and from $z = 0$ to $z = a$.

(b) Repeat for $\mathbf{E} = bx\hat{\mathbf{z}}$.

(c) Repeat for $\mathbf{E} = 13bz\hat{\mathbf{z}} - 7cz\hat{\mathbf{y}}$.

(d) Repeat for $\mathbf{E} = bxz\hat{\mathbf{z}}$.



51 Light is a wave made of electric and magnetic fields, and the fields are perpendicular to the direction of the wave's motion, i.e., they're transverse. An example would be the electric field given by $\mathbf{E} = b\hat{\mathbf{x}} \sin cz$, where b and c are constants. (There would also be an associated magnetic field.) We observe that light can travel through a vacuum, so we expect that this wave pattern is consistent with the nonexistence of any charge in the space it's currently occupying. Use Gauss's law to prove that this is true.



52 This is an alternative approach to problem 49, using a different technique. Suppose that a long cylinder contains a uniform charge density ρ throughout its interior volume.

(a) Use the methods of section 10.7 to find the electric field inside the cylinder. ✓

(b) Extend your solution to the outside region, using the same technique. Once you find the general form of the solution, adjust it so that the inside and outside fields match up at the surface. ✓



53 The purpose of this homework problem is to prove that the divergence is invariant with respect to translations. That is, it

doesn't matter where you choose to put the origin of your coordinate system. Suppose we have a field of the form $\mathbf{E} = ax\hat{\mathbf{x}} + by\hat{\mathbf{y}} + cz\hat{\mathbf{z}}$. This is the most general field we need to consider in any small region as far as the divergence is concerned. (The dependence on x , y , and z is linear, but any smooth function looks linear close up. We also don't need to put in terms like $xy\hat{\mathbf{y}}$, because they don't contribute to the divergence.) Define a new set of coordinates (u, v, w) related to (x, y, z) by

$$\begin{aligned}x &= u + p \\y &= v + q \\z &= w + r,\end{aligned}$$

where p , q , and r are constants. Show that the field's divergence is the same in these new coordinates. Note that $\hat{\mathbf{x}}$ and $\hat{\mathbf{u}}$ are identical, and similarly for the other coordinates. □

54 Using a techniques similar to that of problem 53, show that the divergence is rotationally invariant, in the special case of rotations about the z axis. In such a rotation, we rotate to a new (u, v, z) coordinate system, whose axes are rotated by an angle θ with respect to those of the (x, y, z) system. The coordinates are related by

$$\begin{aligned}x &= u \cos \theta + v \sin \theta \\y &= -u \sin \theta + v \cos \theta\end{aligned}$$

Find how the u and v components the field \mathbf{E} depend on u and v , and show that its divergence is the same in this new coordinate system. □

55 An electric field is given in cylindrical coordinates (R, ϕ, z) by $E_R = ce^{-u|z|}R^{-1}\cos^2\phi$, where the notation E_R indicates the component of the field pointing directly away from the axis, and the components in the other directions are zero. (This isn't a completely impossible expression for the field near a radio transmitting antenna.) (a) Find the total charge enclosed within the infinitely long cylinder extending from the axis out to $R = b$. (b) Interpret the R -dependence of your answer to part a. □

56 Use Euler's theorem to derive the addition theorems that express $\sin(a+b)$ and $\cos(a+b)$ in terms of the sines and cosines of a and b . ▷ Solution, p. 1048 □

57 Find every complex number z such that $z^3 = 1$. ▷ Solution, p. 1048 □

58 Factor the expression $x^3 - y^3$ into factors of the lowest possible order, using complex coefficients. (Hint: use the result of problem 57.) Then do the same using real coefficients. □

59 A dipole consists of two point charges lying on the x axis, a charge $-q$ at the origin, and a $+q$ at $x = \ell$. The dipole is immersed

in an externally imposed, nonuniform electric field with $E_x = bx$, where b is a constant. Add the forces acting on the dipole. Verify that the total force depends only on the dipole moment, not on q or ℓ individually, and that the result is the same as the one found by a fancier method in example 7 on p. 591. \triangleright Solution, p. 1049 ■

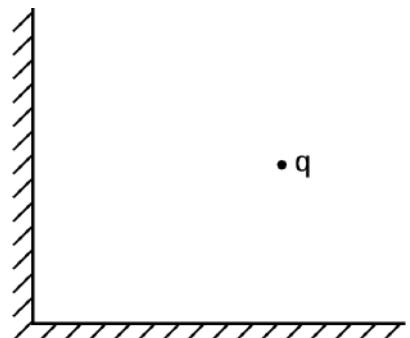
60 In problem 24 on p. 662, you estimated the energy released in a bolt of lightning, based on the energy stored in the electric field immediately before the lightning occurs. The assumption was that the field would build up to a certain value, which is what is necessary to ionize air. However, real-life measurements always seemed to show electric fields strengths roughly 10 times smaller than those required in that model. For a long time, it wasn't clear whether the field measurements were wrong, or the model was wrong. Research carried out in 2003 seems to show that the model was wrong. It is now believed that the final triggering of the bolt of lightning comes from cosmic rays that enter the atmosphere and ionize some of the air. If the field is 10 times smaller than the value assumed in problem 24, what effect does this have on the final result of problem 24? \checkmark ■

61 A charged particle of mass m and charge q is below a horizontal conducting plane. We wish to find the distance ℓ between the particle and the plane so that the particle will be in equilibrium, with its weight supported by electrostatic forces.

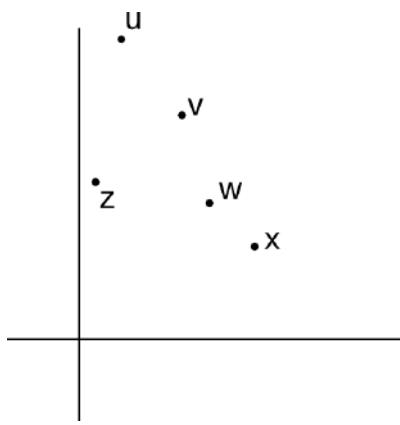
- (a) Determine as much as possible about the form of the answer based on units.
- (b) Find the full result for ℓ .
- (c) Show that the equilibrium is unstable.

62 A point charge q is situated in the empty space inside a corner formed by two perpendicular half-planes made of sheets of metal. Let the sheets lie in the y - z and x - z planes, so that the charge's distances from the planes are x and y . Both x and y are positive. The charge will accelerate due to the electrostatic forces exerted by the sheets. We wish to find the direction θ in which it will accelerate, expressed as an angle counterclockwise from the negative x axis, so that $0 < \theta < \pi/2$.

- (a) Determine as much as possible about the form of the answer based on units.
- (b) Find the full result for θ .



Problem 62.



Problem 63.

63 This problem deals with the cubes and cube roots of complex numbers, but the principles involved apply more generally to other exponents besides 3 and 1/3. These examples are designed to be much easier to do using the magnitude-argument representation of complex numbers than with the cartesian representation. If done by the easiest technique, none of these requires more than two or three lines of *simple* math. In the following, the symbols θ , a , and b represent real numbers, and all angles are to be expressed in radians. As often happens with fractional exponents, the cube root of a complex number will typically have more than one possible value. (Cf. $4^{1/2}$, which can be 2 or -2 .) In parts c and d, this ambiguity is resolved explicitly in the instructions, in a way that is meant to make the calculation as easy as possible.

- (a) Calculate $\arg[(e^{i\theta})^3]$. ✓
- (b) Of the points u , v , w , and x shown in the figure, which could be a cube root of z ?
- (c) Calculate $\arg[\sqrt[3]{a+bi}]$. For simplicity, assume that $a+bi$ is in the first quadrant of the complex plane, and compute the answer for a root that also lies in the first quadrant. ✓
- (d) Compute

$$\frac{1+i}{(-2+2i)^{1/3}}.$$

Because there is more than one possible root to use in the denominator, multiple answers are possible in this problem. Use the root that results in the final answer that lies closest to the real line. (This is also the easiest one to find by using the magnitude-argument techniques introduced in the text.)

✓ ■

- 64** Find the 100th derivative of $e^x \cos x$, evaluated at $x = 0$.
[Based on a problem by T. Needham.] ✓ ■

Key to symbols:

■ easy ■ typical ■ challenging ■ difficult ■ very difficult

✓ An answer check is available at www.lightandmatter.com.

Exercises

Exercise 10A: Field Vectors

Apparatus:

3 solenoids

DC power supply

compass

ruler

cut-off plastic cup

At this point you've studied the gravitational field, \mathbf{g} , and the electric field, \mathbf{E} , but not the magnetic field, \mathbf{B} . However, they all have some of the same mathematical behavior: they act like vectors. Furthermore, magnetic fields are the easiest to manipulate in the lab. Manipulating gravitational fields directly would require futuristic technology capable of moving planet-sized masses around! Playing with electric fields is not as ridiculously difficult, but static electric charges tend to leak off through your body to ground, and static electricity effects are hard to measure numerically. Magnetic fields, on the other hand, are easy to make and control. Any moving charge, i.e., any current, makes a magnetic field.

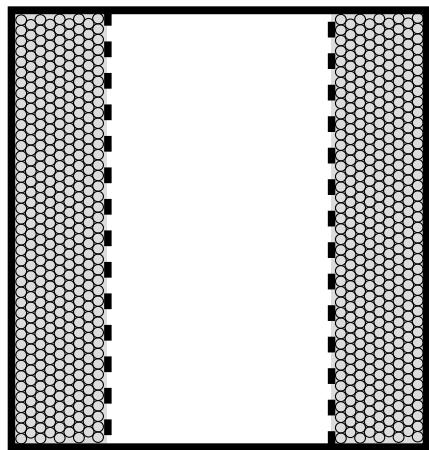
A practical device for making a strong magnetic field is simply a coil of wire, formally known as a solenoid. The field pattern surrounding the solenoid gets stronger or weaker in proportion to the amount of current passing through the wire.

1. With a single solenoid connected to the power supply and laid with its axis horizontal, use a magnetic compass to explore the field pattern inside and outside it. The compass shows you the field vector's direction, but not its magnitude, at any point you choose. Note that the field the compass experiences is a combination (vector sum) of the solenoid's field and the earth's field.
2. What happens when you bring the compass extremely far away from the solenoid?

What does this tell you about the way the solenoid's field varies with distance?

Thus although the compass doesn't tell you the field vector's magnitude numerically, you can get at least some general feel for how it depends on distance.

3. The figure below is a cross-section of the solenoid in the plane containing its axis. Make a sea-of-arrows sketch of the magnetic field in this plane. The length of each arrow should at least approximately reflect the strength of the magnetic field at that point.



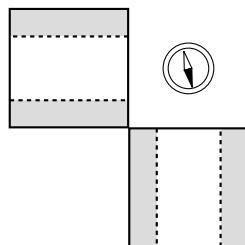
Does the field seem to have sources or sinks?

4. What do you think would happen to your sketch if you reversed the wires?

Try it.

5. Now hook up the two solenoids in parallel. You are going to measure what happens when their two fields combine at a certain point in space. As you've seen already, the solenoids' nearby fields are much stronger than the earth's field; so although we now theoretically have three fields involved (the earth's plus the two solenoids'), it will be safe to ignore the earth's field. The basic idea here is to place the solenoids with their axes at some angle to each other, and put the compass at the intersection of their axes, so that it is the same distance from each solenoid. Since the geometry doesn't favor either solenoid, the only factor that would make one solenoid influence the compass more than the other is current. You can use the cut-off plastic cup as a little platform to bring the compass up to the same level as the solenoids' axes.

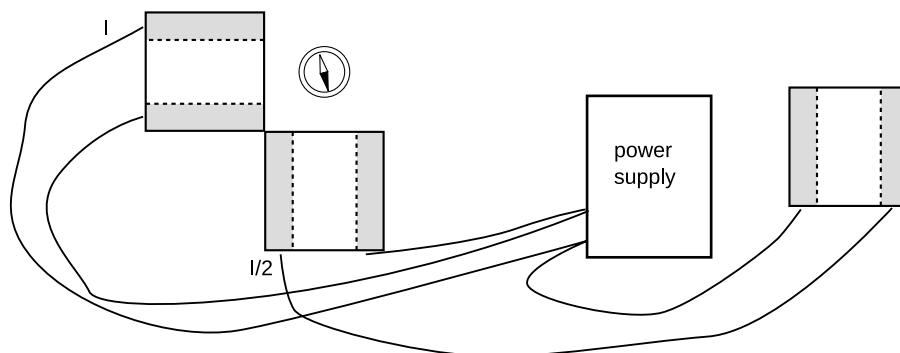
a) What do you think will happen with the solenoids' axes at 90 degrees to each other, and equal currents? Try it. Now represent the vector addition of the two magnetic fields with a diagram. Check your diagram with your instructor to make sure you're on the right track.



b) Now try to make a similar diagram of what would happen if you switched the wires on one of the solenoids.

After predicting what the compass will do, try it and see if you were right.

c) Now suppose you were to go back to the arrangement you had in part a, but you changed one of the currents to half its former value. Make a vector addition diagram, and use trig to predict the angle.



Try it. To cut the current to one of the solenoids in half, an easy and accurate method is simply to put the third solenoid in series with it, and put that third solenoid so far away that its magnetic field doesn't have any significant effect on the compass.

Chapter 11

Electromagnetism

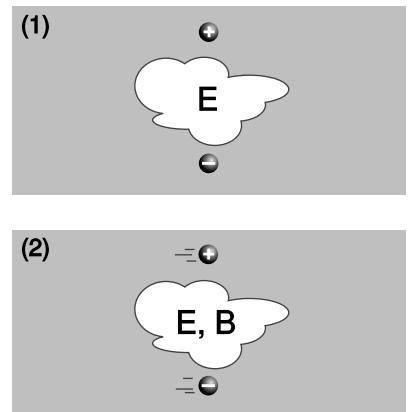
Think not that I am come to destroy the law, or the prophets: I am not come to destroy, but to fulfill.
Matthew 5:17

11.1 More about the magnetic field

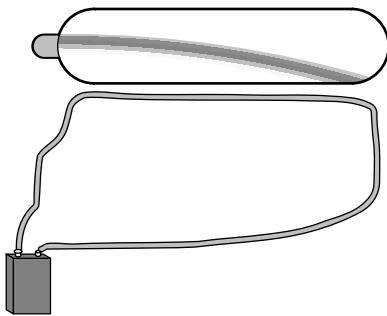
11.1.1 Magnetic forces

In this chapter, I assume you know a few basic ideas about Einstein's theory of relativity, as described in sections 7.1 and 7.2. Unless your typical workday involves rocket ships or particle accelerators, all this relativity stuff might sound like a description of some bizarre futuristic world that is completely hypothetical. There is, however, a relativistic effect that occurs in everyday life, and it is obvious and dramatic: magnetism. Magnetism, as we discussed previously, is an interaction between a moving charge and another moving charge, as opposed to electric forces, which act between any pair of charges, regardless of their motion. Relativistic effects are weak for speeds that are small compared to the speed of light, and the average speed at which electrons drift through a wire is quite low (centimeters per second, typically), so how can relativity be behind an impressive effect like a car being lifted by an electromagnet hanging from a crane? The key is that matter is almost perfectly electrically neutral, and electric forces therefore cancel out almost perfectly. Magnetic forces really aren't very strong, but electric forces are even weaker.

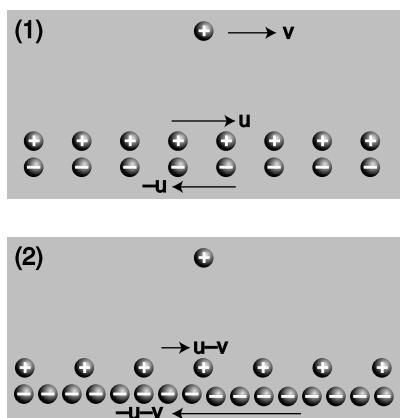
What about the word "relativity" in the name of the theory? It would seem problematic if moving charges interact differently than stationary charges, since motion is a matter of opinion, depending on your frame of reference. Magnetism, however, comes not to destroy relativity but to fulfill it. Magnetic interactions *must* exist according to the theory of relativity. To understand how this can be, consider how time and space behave in relativity. Observers in different frames of reference disagree about the lengths of measuring sticks and the speeds of clocks, but the laws of physics are valid and self-consistent in either frame of reference. Similarly, observers in different frames of reference disagree about what electric and magnetic fields and forces there are, but they agree about concrete physical events. For instance, figure a/1 shows two particles, with opposite charges, which are not moving at a particular mo-



a / The pair of charged particles, as seen in two different frames of reference.



b / A large current is created by shorting across the leads of the battery. The moving charges in the wire attract the moving charges in the electron beam, causing the electrons to curve.



c / A charged particle and a current, seen in two different frames of reference. The second frame is moving at velocity v with respect to the first frame, so all the velocities have v subtracted from them. (As discussed in the main text, this is only approximately correct.)

ment in time. An observer in this frame of reference says there are electric fields around the particles, and predicts that as time goes on, the particles will begin to accelerate towards one another, eventually colliding. A different observer, $a/2$, says the particles are moving. This observer also predicts that the particles will collide, but explains their motion in terms of both an electric field, \mathbf{E} , and a magnetic field, \mathbf{B} . As we'll see shortly, the magnetic field is *required* in order to maintain consistency between the predictions made in the two frames of reference.

To see how this really works out, we need to find a nice simple example that is easy to calculate. An example like figure a is *not* easy to handle, because in the second frame of reference, the moving charges create fields that change over time at any given location. Examples like figure b are easier, because there is a steady flow of charges, and all the fields stay the same over time.¹ What is remarkable about this demonstration is that there can be no electric fields acting on the electron beam at all, since the total charge density throughout the wire is zero. Unlike figure a/2, figure b is purely magnetic.

To see why this must occur based on relativity, we make the mathematically idealized model shown in figure c. The charge by itself is like one of the electrons in the vacuum tube beam of figure b, and a pair of moving, infinitely long line charges has been substituted for the wire. The electrons in a real wire are in rapid thermal motion, and the current is created only by a slow drift superimposed on this chaos. A second deviation from reality is that in the real experiment, the protons are at rest with respect to the tabletop, and it is the electrons that are in motion, but in c/1 we have the positive charges moving in one direction and the negative ones moving the other way. If we wanted to, we could construct a third frame of reference in which the positive charges were at rest, which would be more like the frame of reference fixed to the tabletop in the real demonstration. However, as we'll see shortly, frames c/1 and c/2 are designed so that they are particularly easy to analyze. It's important to note that even though the two line charges are moving in opposite directions, their currents don't cancel. A negative charge moving to the left makes a current that goes to the right, so in frame c/1, the total current is twice that contributed by either line charge.

Frame 1 is easy to analyze because the charge densities of the two line charges cancel out, and the electric field experienced by the

¹For a more practical demonstration of this effect, you can put an ordinary magnet near a computer monitor. The picture will be distorted. Make sure that the monitor has a demagnetizing ("degaussing") button, however! Otherwise you may permanently damage it. Don't use a television tube, because TV tubes don't have demagnetizing buttons.

lone charge is therefore zero:

$$\mathbf{E}_1 = 0$$

In frame 1, any force experienced by the lone charge must therefore be attributed solely to magnetism.

Frame 2 shows what we'd see if we were observing all this from a frame of reference moving along with the lone charge. Why don't the charge densities also cancel in this frame? Here's where the relativity comes in. Relativity tells us that moving objects appear contracted to an observer who is not moving along with them. Both line charges are in motion in both frames of reference, but in frame 1, the line charges were moving at equal speeds, so their contractions were equal, and their charge densities canceled out. In frame 2, however, their speeds are unequal. The positive charges are moving more slowly than in frame 1, so in frame 2 they are less contracted. The negative charges are moving more quickly, so their contraction is greater now. Since the charge densities don't cancel, there is an electric field in frame 2, which points into the wire, attracting the lone charge. Furthermore, the attraction felt by the lone charge must be purely electrical, since the lone charge is at rest in this frame of reference, and magnetic effects occur only between moving charges and other moving charges.²

To summarize, frame 1 displays a purely magnetic attraction, while in frame 2 it is purely electrical.

Now we can calculate the force in frame 2, and equating it to the force in frame 1, we can find out how much magnetic force occurs. To keep the math simple, and to keep from assuming too much about your knowledge of relativity, we're going to carry out this whole calculation in the approximation where all the speeds are fairly small compared to the speed of light.³ For instance, if we find an expression such as $(v/c)^2 + (v/c)^4$, we will assume that the fourth-order term is negligible by comparison. This is known as a calculation "to leading order in v/c ." In fact, I've already used the leading-order approximation twice without saying so! The first

²One could object that this is circular reasoning, since the whole purpose of this argument is to prove from first principles that magnetic effects follow from the theory of relativity. Could there be some extra interaction which occurs between a moving charge and *any* other charge, regardless of whether the other charge is moving or not? We can argue, however, that such a theory would lack self-consistency, since we have to define the electric field somehow, and the only way to define it is in terms of F/q , where F is the force on a test charge q which is at rest. In other words, we'd have to say that there was some extra contribution to the *electric* field if the charge making it was in motion. This would, however, violate Gauss' law, and Gauss' law is amply supported by experiment, even when the sources of the electric field are moving. It would also violate the time-reversal symmetry of the laws of physics.

³The reader who wants to see the full relativistic treatment is referred to E.M. Purcell, *Electricity and Magnetism*, McGraw Hill, 1985, p. 174.

time I used it implicitly was in figure c, where I assumed that the velocities of the two line charges were $u - v$ and $-u - v$. Relativistic velocities don't just combine by simple addition and subtraction like this, but this is an effect we can ignore in the present approximation. The second sleight of hand occurred when I stated that we could equate the forces in the two frames of reference. Force, like time and distance, is distorted relativistically when we change from one frame of reference to another. Again, however, this is an effect that we can ignore to the desired level of approximation.

Let $\pm\lambda$ be the charge per unit length of each line charge without relativistic contraction, i.e., in the frame moving with that line charge. Using the approximation $\gamma = (1 - v^2/c^2)^{-1/2} \approx 1 + v^2/2c^2$ for $v \ll c$, the total charge per unit length in frame 2 is

$$\begin{aligned}\lambda_{total, 2} &\approx \lambda \left[1 + \frac{(u - v)^2}{2c^2} \right] - \lambda \left[1 + \frac{(-u - v)^2}{2c^2} \right] \\ &= \frac{-2\lambda uv}{c^2}.\end{aligned}$$

Let R be the distance from the line charge to the lone charge. Applying Gauss' law to a cylinder of radius R centered on the line charge, we find that the magnitude of the electric field experienced by the lone charge in frame 2 is

$$E = \frac{4k\lambda uv}{c^2 R},$$

and the force acting on the lone charge q is

$$F = \frac{4k\lambda quv}{c^2 R}.$$

In frame 1, the current is $I = 2\lambda_1 u$ (see homework problem 5), which we can approximate as $I = 2\lambda u$, since the current, unlike $\lambda_{total, 2}$, doesn't vanish completely without the relativistic effect. The magnetic force on the lone charge q due to the current I is

$$F = \frac{2kIqv}{c^2 R}.$$

Discussion Questions

A In the situation shown in figure c, is there a frame in which the force \mathbf{F} is a purely electric one, \mathbf{F}_E ? Pure \mathbf{F}_B ?

Is there a frame in which the electromagnetic field is a pure \mathbf{E} ? Pure \mathbf{B} ?

Is the charge density ρ zero in both frames? One? Neither?

What about the current I (or current density \mathbf{j})?

B For the situation shown in figure c, draw a spacetime diagram in the style demonstrated in sec. 7.2, p. 400, showing the positive charges as black world-lines and the negative as red, in the wire's rest frame. Use a ruler, and draw the spacing fairly accurately. Interpret this in the frame of the lone charge.

C Resolve the following paradox concerning the argument given in this section. We would expect that at any given time, electrons in a solid would be associated with protons in a definite way. For simplicity, let's imagine that the solid is made out of hydrogen (which actually does become a metal under conditions of very high pressure). A hydrogen atom consists of a single proton and a single electron. Even if the electrons are moving and forming an electric current, we would imagine that this would be like a game of musical chairs, with the protons as chairs and the electrons as people. Each electron has a proton that is its “friend,” at least for the moment. This is the situation shown in figure c/1. How, then, can an observer in a different frame see the electrons and protons as not being paired up, as in c/2?

11.1.2 The magnetic field

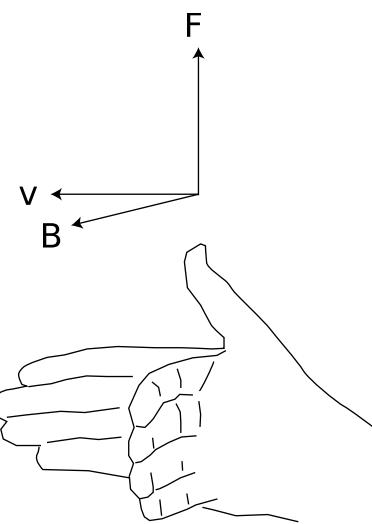
Definition in terms of the force on a moving particle

With electricity, it turned out to be useful to define an electric field rather than always working in terms of electric forces. Likewise, we want to define a magnetic field, **B**. Let's look at the result of the preceding subsection for insight. The equation

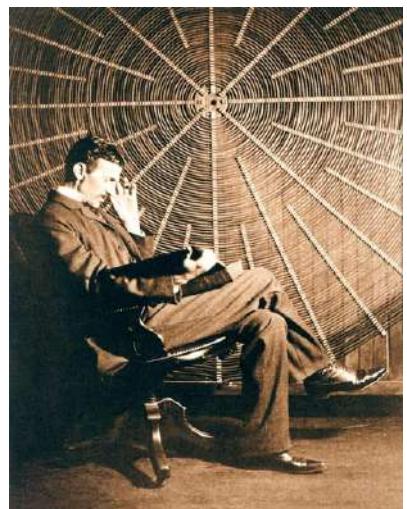
$$F = \frac{2kIqv}{c^2R}$$

shows that when we put a moving charge near other moving charges, there is an extra magnetic force on it, in addition to any electric forces that may exist. Equations for electric forces always have a factor of k in front — the Coulomb constant k is called the coupling constant for electric forces. Since magnetic effects are relativistic in origin, they end up having a factor of k/c^2 instead of just k . In a world where the speed of light was infinite, relativistic effects, including magnetism, would be absent, and the coupling constant for magnetism would be zero. A cute feature of the metric system is that we have $k/c^2 = 10^{-7} \text{ N}\cdot\text{s}^2/\text{C}^2$ exactly, as a matter of definition.

Naively, we could try to work by analogy with the electric field, and define the magnetic field as the magnetic force per unit charge. However, if we think of the lone charge in our example as the test charge, we'll find that this approach fails, because the force depends not just on the test particle's charge, but on its velocity, v , as well. Although we only carried out calculations for the case where the particle was moving parallel to the wire, in general this velocity is a vector, \mathbf{v} , in three dimensions. We can also anticipate that the magnetic field will be a vector. The electric and gravitational fields are vectors, and we expect intuitively based on our experience with magnetic compasses that a magnetic field has a particular direction in space. Furthermore, reversing the current I in our example would have reversed the force, which would only make sense if the magnetic field had a direction in space that could be reversed. Summarizing, we think there must be a magnetic field vector **B**, and the force on a test particle moving through a magnetic field is proportional both to the **B** vector and to the particle's own **v** vector. In other



d / The right-hand relationship between the velocity of a positively charged particle, the magnetic field through which it is moving, and the magnetic force on it.

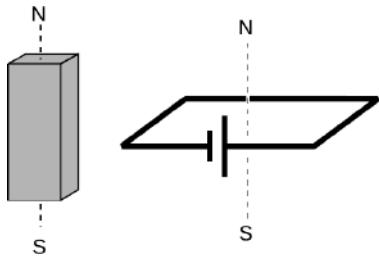


e / The unit of magnetic field, the tesla, is named after Serbian-American inventor Nikola Tesla.

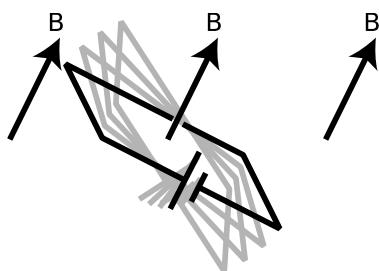
words, the magnetic force vector \mathbf{F} is found by some sort of vector multiplication of the vectors \mathbf{v} and \mathbf{B} . As proved on page 1027, however, there is only one physically useful way of defining such a multiplication, which is the cross product.

We therefore define the magnetic field vector, \mathbf{B} , as the vector that determines the force on a charged particle according to the following rule:

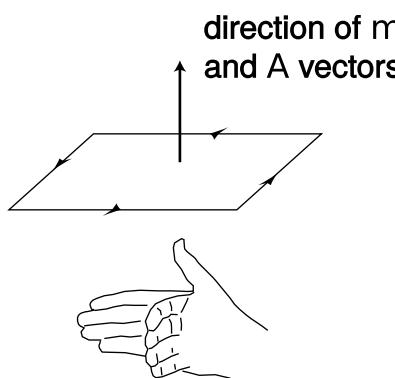
$$\mathbf{F} = q\mathbf{v} \times \mathbf{B} \quad [\text{definition of the magnetic field}]$$



f / A standard dipole made from a square loop of wire shorting across a battery. It acts very much like a bar magnet, but its strength is more easily quantified.



g / A dipole tends to align itself to the surrounding magnetic field.



h / The \mathbf{m} and \mathbf{A} vectors.

From this definition, we see that the magnetic field's units are $\text{N} \cdot \text{s/C} \cdot \text{m}$, which are usually abbreviated as teslas, $1 \text{ T} = 1 \text{ N} \cdot \text{s/C} \cdot \text{m}$. The definition implies a right-hand-rule relationship among the vectors, figure d, if the charge q is positive, and the opposite handedness if it is negative.

This is not just a definition but a bold prediction! Is it really true that for any point in space, we can always find a vector \mathbf{B} that successfully predicts the force on any passing particle, regardless of its charge and velocity vector? Yes — it's not obvious that it can be done, but experiments verify that it can. How? Well for example, the cross product of parallel vectors is zero, so we can try particles moving in various directions, and hunt for the direction that produces zero force; the \mathbf{B} vector lies along that line, in either the same direction the particle was moving, or the opposite one. We can then go back to our data from one of the other cases, where the force was nonzero, and use it to choose between these two directions and find the magnitude of the \mathbf{B} vector. We could then verify that this vector gave correct force predictions in a variety of other cases.

Even with this empirical reassurance, the meaning of this equation is not intuitively transparent, nor is it practical in most cases to measure a magnetic field this way. For these reasons, let's look at an alternative method of defining the magnetic field which, although not as fundamental or mathematically simple, may be more appealing.

Definition in terms of the torque on a dipole

A compass needle in a magnetic field experiences a torque which tends to align it with the field. This is just like the behavior of an electric dipole in an electric field, so we consider the compass needle to be a *magnetic dipole*. In subsection 10.1.3 on page 590, we gave an alternative definition of the electric field in terms of the torque on an electric dipole.

To define the strength of a magnetic field, however, we need some way of defining the strength of a test dipole, i.e., we need a definition of the magnetic dipole moment. We could use an iron permanent magnet constructed according to certain specifications,

but such an object is really an extremely complex system consisting of many iron atoms, only some of which are aligned with each other. A more fundamental standard dipole is a square current loop. This could be little resistive circuit consisting of a square of wire shorting across a battery, f.

Applying $\mathbf{F} = \mathbf{v} \times \mathbf{B}$, we find that such a loop, when placed in a magnetic field, g, experiences a torque that tends to align plane so that its interior “face” points in a certain direction. Since the loop is symmetric, it doesn’t care if we rotate it like a wheel without changing the plane in which it lies. It is this preferred facing direction that we will end up using as our alternative definition of the magnetic field.

If the loop is out of alignment with the field, the torque on it is proportional to the amount of current, and also to the interior area of the loop. The proportionality to current makes sense, since magnetic forces are interactions between moving charges, and current is a measure of the motion of charge. The proportionality to the loop’s area is also not hard to understand, because increasing the length of the sides of the square increases both the amount of charge contained in this circular “river” and the amount of leverage supplied for making torque. Two separate physical reasons for a proportionality to length result in an overall proportionality to length squared, which is the same as the area of the loop. For these reasons, we define the magnetic dipole moment of a square current loop as

$$\mathbf{m} = IA,$$

where the direction of the vectors is defined as shown in figure h.

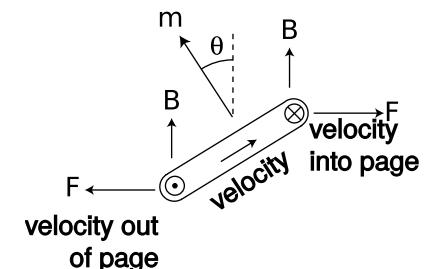
We can now give an alternative definition of the magnetic field:

The magnetic field vector, \mathbf{B} , at any location in space is defined by observing the torque exerted on a magnetic test dipole \mathbf{m}_t consisting of a square current loop. The field’s magnitude is

$$|\mathbf{B}| = \frac{\tau}{|\mathbf{m}_t| \sin \theta},$$

where θ is the angle between the dipole vector and the field. This is equivalent to the vector cross product $\tau = \mathbf{m}_t \times \mathbf{B}$.

Let’s show that this is consistent with the previous definition, using the geometry shown in figure i. The velocity vectors that point in and out of the page are shown using the convention defined in figure j. Let the mobile charge carriers in the wire have linear density λ , and let the sides of the loop have length h , so that we have $I = \lambda v$, and $m = h^2 \lambda v$. The only nonvanishing torque comes from the forces on the left and right sides. The currents in these sides are perpendicular to the field, so the magnitude of the cross

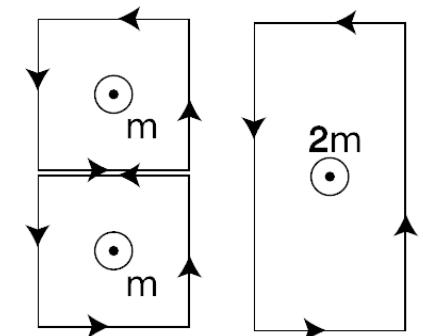


i / The torque on a current loop in a magnetic field. The current comes out of the page, goes across, goes back into the page, and then back across the other way in the hidden side of the loop.

⊖ out of the page

⊗ into the page

j / A vector coming out of the page is shown with the tip of an arrowhead. A vector going into the page is represented using the tailfeathers of the arrow.



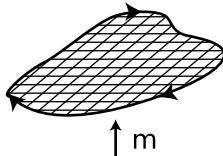
k / Dipole vectors can be added.

product $\mathbf{F} = q\mathbf{v} \times \mathbf{B}$ is simply $|\mathbf{F}| = qvB$. The torque supplied by each of these forces is $\mathbf{r} \times \mathbf{F}$, where the lever arm \mathbf{r} has length $h/2$, and makes an angle θ with respect to the force vector. The magnitude of the total torque acting on the loop is therefore

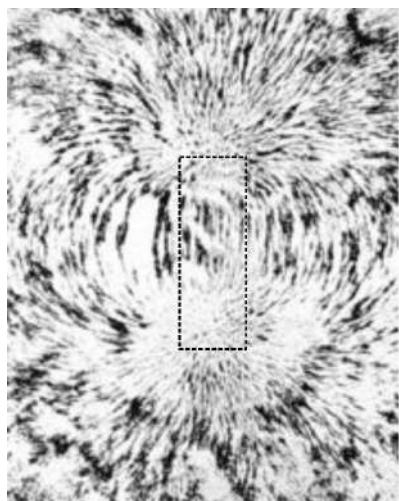
$$|\tau| = 2 \frac{h}{2} |\mathbf{F}| \sin \theta \\ = h qvB \sin \theta,$$

and substituting $q = \lambda h$ and $v = m/h^2\lambda$, we have

$$|\tau| = h \lambda h \frac{m}{h^2 \lambda} B \sin \theta \\ = mB \sin \theta,$$



l / An irregular loop can be broken up into little squares.



m / The magnetic field pattern around a bar magnet is created by the superposition of the dipole fields of the individual iron atoms. Roughly speaking, it looks like the field of one big dipole, especially farther away from the magnet. Closer in, however, you can see a hint of the magnet's rectangular shape. The picture was made by placing iron filings on a piece of paper, and then bringing a magnet up underneath.

which is consistent with the second definition of the field.

It undoubtedly seems artificial to you that we have discussed dipoles only in the form of a square loop of current. A permanent magnet, for example, is made out of atomic dipoles, and atoms aren't square! However, it turns out that the shape doesn't matter. To see why this is so, consider the additive property of areas and dipole moments, shown in figure k. Each of the square dipoles has a dipole moment that points out of the page. When they are placed side by side, the currents in the adjoining sides cancel out, so they are equivalent to a single rectangular loop with twice the area. We can break down any irregular shape into little squares, as shown in figure l, so the dipole moment of any planar current loop can be calculated based on its area, regardless of its shape.

The magnetic dipole moment of an atom

example 1

Let's make an order-of-magnitude estimate of the magnetic dipole moment of an atom. A hydrogen atom is about 10^{-10} m in diameter, and the electron moves at speeds of about 10^{-2} c. We don't know the shape of the orbit, and indeed it turns out that according to the principles of quantum mechanics, the electron doesn't even have a well-defined orbit, but if we're brave, we can still estimate the dipole moment using the cross-sectional area of the atom, which will be on the order of $(10^{-10} \text{ m})^2 = 10^{-20} \text{ m}^2$. The electron is a single particle, not a steady current, but again we throw caution to the winds, and estimate the current it creates as $e/\Delta t$, where Δt , the time for one orbit, can be estimated by dividing the size of the atom by the electron's velocity. (This is only a rough estimate, and we don't know the shape of the orbit, so it would be silly, for instance, to bother with multiplying the diameter by π based on our intuitive visualization of the electron as moving around the circumference of a circle.) The result for the dipole moment is $m \sim 10^{-23} \text{ A}\cdot\text{m}^2$.

Should we be impressed with how small this dipole moment is, or with how big it is, considering that it's being made by a single atom? Very large or very small numbers are never very interesting by themselves. To get a feeling for what they mean, we need to compare them to something else. An interesting comparison here is to think in terms of the total number of atoms in a typical object, which might be on the order of 10^{26} (Avogadro's number). Suppose we had this many atoms, with their moments all aligned. The total dipole moment would be on the order of $10^3 \text{ A}\cdot\text{m}^2$, which is a pretty big number. To get a dipole moment this strong using human-scale devices, we'd have to send a thousand amps of current through a one-square meter loop of wire! The insight to be gained here is that, even in a permanent magnet, we must not have all the atoms perfectly aligned, because that would cause more spectacular magnetic effects than we really observe. Apparently, nearly all the atoms in such a magnet are oriented randomly, and do not contribute to the magnet's dipole moment.

Discussion Questions

A The physical situation shown in figure c on page 676 was analyzed entirely in terms of forces. Now let's go back and think about it in terms of fields. The charge by itself up above the wire is like a test charge, being used to determine the magnetic and electric fields created by the wire. In figures c/1 and c/2, are there fields that are purely electric or purely magnetic? Are there fields that are a mixture of **E** and **B**? How does this compare with the forces?

B Continuing the analysis begun in discussion question A, can we come up with a scenario involving some charged particles such that the fields are purely magnetic in one frame of reference but a mixture of **E** and **B** in another frame? How about an example where the fields are purely electric in one frame, but mixed in another? Or an example where the fields are purely electric in one frame, but purely magnetic in another?

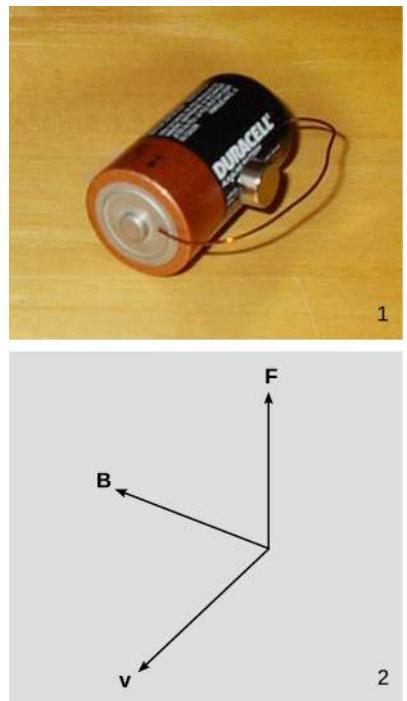
11.1.3 Some applications

Magnetic levitation

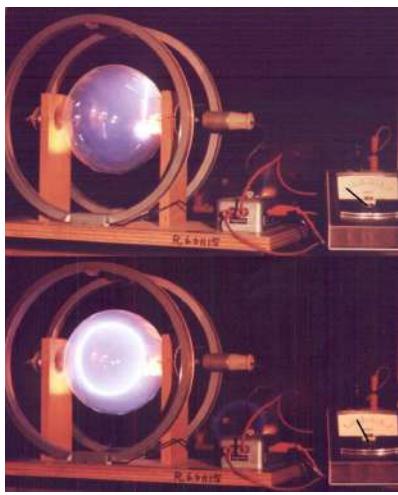
example 2

In figure n, a small, disk-shaped permanent magnet is stuck on the side of a battery, and a wire is clasped loosely around the battery, shorting it. A large current flows through the wire. The electrons moving through the wire feel a force from the magnetic field made by the permanent magnet, and this force levitates the wire.

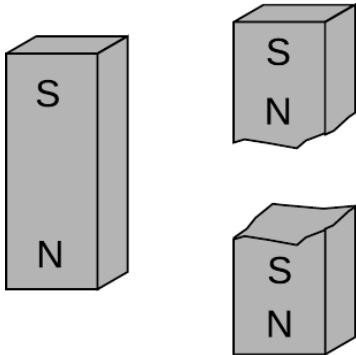
From the photo, it's possible to find the direction of the magnetic field made by the permanent magnet. The electrons in the copper wire are negatively charged, so they flow from the negative (flat) terminal of the battery to the positive terminal (the one with the bump, in front). As the electrons pass by the permanent magnet, we can imagine that they would experience a field either toward the magnet, or away from it, depending on which way the magnet



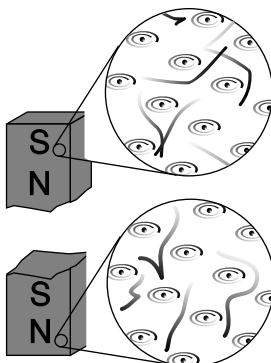
n / Example 2.



o / Magnetic forces cause a beam of electrons to move in a circle.



p / You can't isolate the poles of a magnet by breaking it in half.



q / A magnetic dipole is made out of other dipoles, not out of monopoles.

was flipped when it was stuck onto the battery. By the right-hand rule (figure d on page 679), the field must be toward the battery.

Nervous-system effects during an MRI scan example 3

During an MRI scan of the head, the patient's nervous system is exposed to intense magnetic fields, and there are ions moving around in the nerves. The resulting forces on the ions can cause symptoms such as vertigo.

A circular orbit example 4

The magnetic force is always perpendicular to the motion of the particle, so it can never do any work, and a charged particle moving through a magnetic field does not experience any change in its kinetic energy: its velocity vector can change its direction, but not its magnitude. If the velocity vector is initially perpendicular to the field, then the curve of its motion will remain in the plane perpendicular to the field, so the magnitude of the magnetic force on it will stay the same. When an object experiences a force with constant magnitude, which is always perpendicular to the direction of its motion, the result is that it travels in a circle.

Figure o shows a beam of electrons in a spherical vacuum tube. In the top photo, the beam is emitted near the right side of the tube, and travels straight up. In the bottom photo, a magnetic field has been imposed by an electromagnet surrounding the vacuum tube; the ammeter on the right shows that the current through the electromagnet is now nonzero. We observe that the beam is bent into a circle.

self-check A

Infer the direction of the magnetic field. Don't forget that the beam is made of electrons, which are negatively charged! ➤ Answer, p. 1064

Homework problem 12 is a quantitative analysis of circular orbits.

A velocity filter example 5

Suppose you see the electron beam in figure o, and you want to determine how fast the electrons are going. You certainly can't do it with a stopwatch! Physicists may also encounter situations where they have a beam of unknown charged particles, and they don't even know their charges. This happened, for instance, when alpha and beta radiation were discovered. One solution to this problem relies on the fact that the force experienced by a charged particle in an electric field, $\mathbf{F}_E = q\mathbf{E}$, is independent of its velocity, but the force due to a magnetic field, $\mathbf{F}_B = q\mathbf{v} \times \mathbf{B}$, isn't. One can send a beam of charged particles through a space containing both an electric and a magnetic field, setting up the fields so that the two forces will cancel out perfectly for a certain velocity. Note that since both forces are proportional to the charge of the particles, the cancellation is independent of charge. Such a *velocity filter* can be used either to determine the velocity of an unknown

beam or particles, or to select from a beam of particles only those having velocities within a certain desired range. Homework problem 7 is an analysis of this application.

11.1.4 No magnetic monopoles

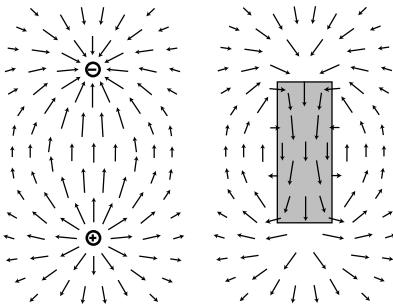
If you could play with a handful of electric dipoles and a handful of bar magnets, they would appear very similar. For instance, a pair of bar magnets wants to align themselves head-to-tail, and a pair of electric dipoles does the same thing. (It is unfortunately not that easy to make a permanent electric dipole that can be handled like this, since the charge tends to leak.)

You would eventually notice an important difference between the two types of objects, however. The electric dipoles can be broken apart to form isolated positive charges and negative charges. The two-ended device can be broken into parts that are not two-ended. But if you break a bar magnet in half, p, you will find that you have simply made two smaller two-ended objects.

The reason for this behavior is not hard to divine from our microscopic picture of permanent iron magnets. An electric dipole has extra positive “stuff” concentrated in one end and extra negative in the other. The bar magnet, on the other hand, gets its magnetic properties not from an imbalance of magnetic “stuff” at the two ends but from the orientation of the rotation of its electrons. One end is the one from which we could look down the axis and see the electrons rotating clockwise, and the other is the one from which they would appear to go counterclockwise. There is no difference between the “stuff” in one end of the magnet and the other, q.

Nobody has ever succeeded in isolating a single magnetic pole. In technical language, we say that magnetic *monopoles* not seem to exist. Electric monopoles *do* exist — that’s what charges are.

Electric and magnetic forces seem similar in many ways. Both act at a distance, both can be either attractive or repulsive, and both are intimately related to the property of matter called charge. (Recall that magnetism is an interaction between moving charges.) Physicists’s aesthetic senses have been offended for a long time because this seeming symmetry is broken by the existence of electric monopoles and the absence of magnetic ones. Perhaps some exotic form of matter exists, composed of particles that are magnetic monopoles. If such particles could be found in cosmic rays or moon rocks, it would be evidence that the apparent asymmetry was only an asymmetry in the composition of the universe, not in the laws of physics. For these admittedly subjective reasons, there have been several searches for magnetic monopoles. Experiments have been performed, with negative results, to look for magnetic monopoles embedded in ordinary matter. Soviet physicists in the 1960’s made exciting claims that they had created and detected mag-



r / Magnetic fields have no sources or sinks.

netic monopoles in particle accelerators, but there was no success in attempts to reproduce the results there or at other accelerators. The most recent search for magnetic monopoles, done by reanalyzing data from the search for the top quark at Fermilab, turned up no candidates, which shows that either monopoles don't exist in nature or they are extremely massive and thus hard to create in accelerators.

The nonexistence of magnetic monopoles means that unlike an electric field, a magnetic one, can never have sources or sinks. The magnetic field vectors lead in paths that loop back on themselves, without ever converging or diverging at a point, as in the fields shown in figure r. Gauss' law for magnetism is therefore much simpler than Gauss' law for electric fields:

$$\Phi_B = \sum \mathbf{B}_j \cdot \mathbf{A}_j = 0$$

The magnetic flux through any closed surface is zero.

self-check B

Draw a Gaussian surface on the electric dipole field of figure r that has nonzero electric flux through it, and then draw a similar surface on the magnetic field pattern. What happens? ▷ Answer, p. 1064

The field of a wire

example 6

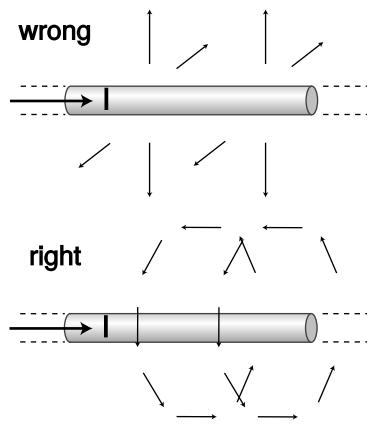
▷ On page 678, we showed that a long, straight wire carrying current I exerts a magnetic force

$$F = \frac{2kIqv}{c^2R}$$

on a particle with charge q moving parallel to the wire with velocity v . What, then, is the magnetic field of the wire?

▷ Comparing the equation above to the first definition of the magnetic field, $\mathbf{F} = \mathbf{v} \times \mathbf{B}$, it appears that the magnetic field is one that falls off like $1/R$, where R is the distance from the wire. However, it's not so easy to determine the direction of the field vector. There are two other axes along which the particle could have been moving, and the brute-force method would be to carry out relativistic calculations for these cases as well. Although this would probably be enough information to determine the field, we don't want to do that much work.

Instead, let's consider what the possibilities are. The field can't be parallel to the wire, because a cross product vanishes when the two vectors are parallel, and yet we know from the case we analyzed that the force doesn't vanish when the particle is moving parallel to the wire. The other two possibilities that are consistent with the symmetry of the problem are shown in figure s. One is like a bottle brush, and the other is like a spool of thread. The bottle brush pattern, however, violates Gauss' law for magnetism. If we made a cylindrical Gaussian surface with its axis coinciding



s / Example 6.

with the wire, the flux through it would *not* be zero. We therefore conclude that the spool-of-thread pattern is the correct one.⁴ Since the particle in our example was moving perpendicular to the field, we have $|F| = |q||v||B|$, so

$$\begin{aligned}|B| &= \frac{|F|}{|q||v|} \\ &= \frac{2kI}{c^2R}\end{aligned}$$

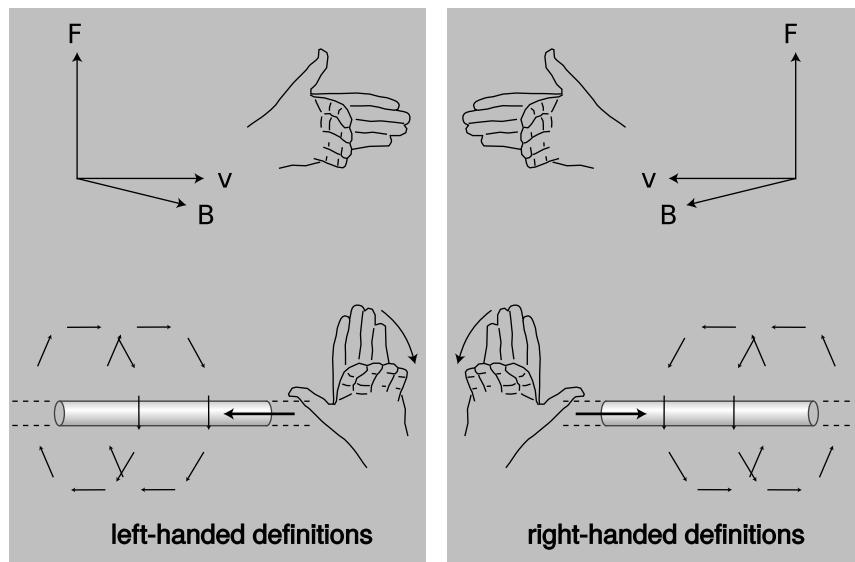
11.1.5 Symmetry and handedness

Imagine that you establish radio contact with an alien on another planet. Neither of you even knows where the other one's planet is, and you aren't able to establish any landmarks that you both recognize. You manage to learn quite a bit of each other's languages, but you're stumped when you try to establish the definitions of left and right (or, equivalently, clockwise and counterclockwise). Is there any way to do it?

If there was any way to do it without reference to external landmarks, then it would imply that the laws of physics themselves were asymmetric, which would be strange. Why should they distinguish left from right? The gravitational field pattern surrounding a star or planet looks the same in a mirror, and the same goes for electric fields. However, the magnetic field patterns shown in figure s seems to violate this principle. Could you use these patterns to explain left and right to the alien? No. If you look back at the definition of the magnetic field, it also contains a reference to handedness: the direction of the vector cross product. The aliens might have reversed their definition of the magnetic field, in which case their drawings of field patterns would look like mirror images of ours, as in the left panel of figure t.

Until the middle of the twentieth century, physicists assumed that any reasonable set of physical laws would have to have this kind of symmetry between left and right. An asymmetry would

⁴Strictly speaking, there is a hole in this logic, since I've only ruled out a field that is purely along one of these three perpendicular directions. What if it has components along more than one of them? A little more work is required to eliminate these mixed possibilities. For example, we can rule out a field with a nonzero component parallel to the wire based on the following symmetry argument. Suppose a charged particle is moving in the plane of the page directly toward the wire. If the field had a component parallel to the wire, then the particle would feel a force into or out of the page, but such a force is impossible based on symmetry, since the whole arrangement is symmetric with respect to mirror-reflection across the plane of the page.



t / Left-handed and right-handed definitions.



u / In this scene from Swan Lake, the choreography has a symmetry with respect to left and right.



v / C.S. Wu

be grotesque. Whatever their aesthetic feelings, they had to change their opinions about reality when experiments by C.S. Wu et al. showed that the weak nuclear force violates right-left symmetry! It is still a mystery why right-left symmetry is observed so scrupulously in general, but is violated by one particular type of physical process.

11.2 Magnetic fields by superposition

11.2.1 Superposition of straight wires

In chapter 10, one of the most important goals was to learn how to calculate the electric field for a given charge distribution. The corresponding problem for magnetism would be to calculate the magnetic field arising from a given set of currents. So far, however, we only know how to calculate the magnetic field of a long, straight wire,

$$B = \frac{2kI}{c^2 R},$$

with the geometry shown in figure a. Whereas a charge distribution can be broken down into individual point charges, most currents cannot be broken down into a set of straight-line currents. Nevertheless, let's see what we can do with the tools that we have.

A ground fault interrupter

example 7

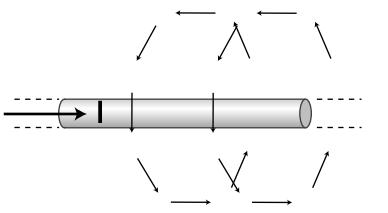
Electric current in your home is supposed to flow out of one side of the outlet, through an appliance, and back into the wall through the other side of the outlet. If that's not what happens, then we have a problem — the current must be finding its way to ground through some other path, perhaps through someone's body. If you have outlets in your home that have "test" and "reset" buttons on them, they have a safety device built into them that is meant to protect you in this situation. The ground fault interrupter (GFI) shown in figure b, routes the outgoing and returning currents through two wires that lie very close together. The clockwise and counterclockwise fields created by the two wires combine by vector addition, and normally cancel out almost exactly. However, if current is not coming back through the circuit, a magnetic field is produced. The doughnut-shaped collar detects this field (using an effect called induction, to be discussed in section 11.5), and sends a signal to a logic chip, which breaks the circuit within about 25 milliseconds.

An example with vector addition

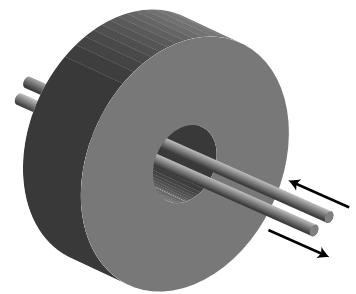
example 8

- ▷ Two long, straight wires each carry current I parallel to the y axis, but in opposite directions. They are separated by a gap $2h$ in the x direction. Find the magnitude and direction of the magnetic field at a point located at a height z above the plane of the wires, directly above the center line.
- ▷ The magnetic fields contributed by the two wires add like vectors, which means we can add their x and z components. The x components cancel by symmetry. The magnitudes of the individual fields are equal,

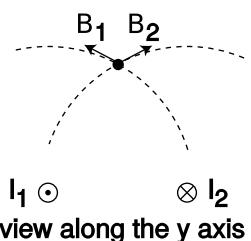
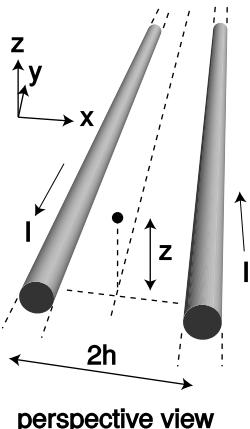
$$B_1 = B_2 = \frac{2kI}{c^2 R},$$



a / The magnetic field of a long, straight wire.



b / A ground fault interrupter.



c / Example 8.

so the total field in the z direction is

$$B_z = 2 \frac{2kI}{c^2 R} \sin \theta,$$

where θ is the angle the field vectors make above the x axis. The sine of this angle equals h/R , so

$$B_z = \frac{4kIh}{c^2 R^2}.$$

(Putting this explicitly in terms of z gives the less attractive form $B_z = 4kIh/c^2(h^2 + z^2)$.)

At large distances from the wires, the individual fields are mostly in the $\pm x$ direction, so most of their strength cancels out. It's not surprising that the fields tend to cancel, since the currents are in opposite directions. What's more interesting is that not only is the field weaker than the field of one wire, it also falls off as R^{-2} rather than R^{-1} . If the wires were right on top of each other, their currents would cancel each other out, and the field would be zero. From far away, the wires appear to be almost on top of each other, which is what leads to the more drastic R^{-2} dependence on distance.

self-check C

In example 8, what is the field right between the wires, at $z = 0$, and how does this simpler result follow from vector addition? \triangleright Answer, p. 1064

An alarming infinity

An interesting aspect of the R^{-2} dependence of the field in example 8 is the energy of the field. We've already established on p. 612 that the energy density of the magnetic field must be proportional to the square of the field strength, B^2 , the same as for the gravitational and electric fields. Suppose we try to calculate the energy per unit length stored in the field of a *single* wire. We haven't yet found the proportionality factor that goes in front of the B^2 , but that doesn't matter, because the energy per unit length turns out to be infinite! To see this, we can construct concentric cylindrical shells of length L , with each shell extending from R to $R + dR$. The volume of the shell equals its circumference times its thickness times its length, $dv = (2\pi R)(dR)(L) = 2\pi L dR$. For a single wire, we have $B \sim R^{-1}$, so the energy density is proportional to R^{-2} , and the energy contained in each shell varies as $R^{-2} dv \sim R^{-1} dr$. Integrating this gives a logarithm, and as we let R approach infinity, we get the logarithm of infinity, which is infinite.

Taken at face value, this result would imply that electrical currents could never exist, since establishing one would require an infinite amount of energy per unit length! In reality, however, we would be dealing with an electric *circuit*, which would be more like

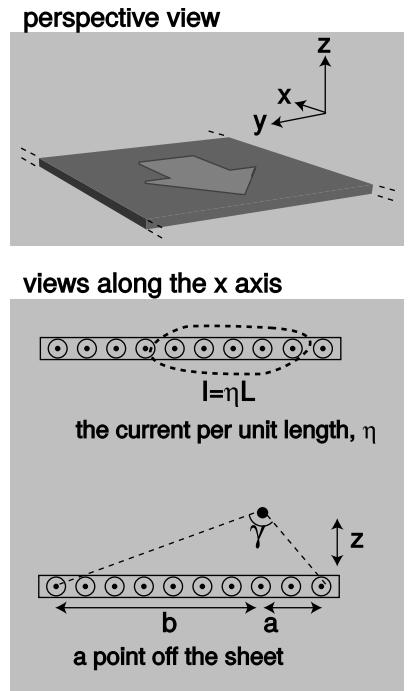
the two wires of example 8: current goes out one wire, but comes back through the other. Since the field really falls off as R^{-2} , we have an energy density that varies as R^{-4} , which does *not* give infinity when integrated out to infinity. (There is still an infinity at $R = 0$, but this doesn't occur for a real wire, which has a finite diameter.)

Still, one might worry about the physical implications of the single-wire result. For instance, suppose we turn on an electron gun, like the one in a TV tube. It takes perhaps a microsecond for the beam to progress across the tube. After it hits the other side of the tube, a return current is established, but at least for the first microsecond, we have only a single current, not two. Do we have infinite energy in the resulting magnetic field? No. It takes time for electric and magnetic field disturbances to travel outward through space, so during that microsecond, the field spreads only to some finite value of R , not $R = \infty$.

This reminds us of an important fact about our study of magnetism so far: we have only been considering situations where the currents and magnetic fields are constant over time. The equation $B = 2kI/c^2R$ was derived under this assumption. This equation is only valid if we assume the current has been established and flowing steadily for a long time, and if we are talking about the field at a point in space at which the field has been established for a long time. The generalization to time-varying fields is nontrivial, and qualitatively new effects will crop up. We have already seen one example of this on page 622, where we inferred that an inductor's time-varying magnetic field creates an electric field — an electric field which is not created by any charges anywhere. Effects like these will be discussed in section 11.5.

A sheet of current

There is a saying that in computer science, there are only three nice numbers: zero, one, and however many you please. In other words, computer software shouldn't have arbitrary limitations like a maximum of 16 open files, or 256 e-mail messages per mailbox. When superposing the fields of long, straight wires, the really interesting cases are one wire, two wires, and infinitely many wires. With an infinite number of wires, each carrying an infinitesimal current, we can create sheets of current, as in figure d. Such a sheet has a certain amount of current per unit length, which we notate η (Greek letter eta). The setup is similar to example 8, except that all the currents are in the same direction, and instead of adding up two fields, we add up an infinite number of them by doing an integral.

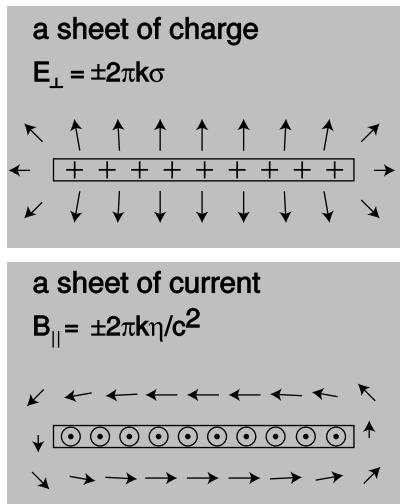


d / A sheet of charge.

For the y component, we have

$$\begin{aligned}
 B_y &= \int \frac{2k dI}{c^2 R} \cos \theta \\
 &= \int_{-a}^b \frac{2k\eta dy}{c^2 R} \cos \theta \\
 &= \frac{2k\eta}{c^2} \int_{-a}^b \frac{\cos \theta}{R} dy \\
 &= \frac{2k\eta}{c^2} \int_{-a}^b \frac{z dy}{y^2 + z^2} \\
 &= \frac{2k\eta}{c^2} \left(\tan^{-1} \frac{b}{z} - \tan^{-1} \frac{-a}{z} \right) \\
 &= \frac{2k\eta\gamma}{c^2},
 \end{aligned}$$

where in the last step we have used the identity $\tan^{-1}(-x) = -\tan^{-1}x$, combined with the relation $\tan^{-1}b/z + \tan^{-1}a/z = \gamma$, which can be verified with a little geometry and trigonometry. The calculation of B_z is left as an exercise (problem 23). More interesting is what happens underneath the sheet: by the right-hand rule, all the currents make rightward contributions to the field there, so B_y abruptly reverses itself as we pass through the sheet.



e / A sheet of charge and a sheet of current.

Close to the sheet, the angle γ approaches π , so we have

$$B_y = \frac{2\pi k\eta}{c^2}.$$

Figure e shows the similarity between this result and the result for a sheet of charge. In one case the sources are charges and the field is electric; in the other case we have currents and magnetic fields. In both cases we find that the field changes suddenly when we pass through a sheet of sources, and the amount of this change doesn't depend on the size of the sheet. It was this type of reasoning that eventually led us to Gauss' law in the case of electricity, and in section 11.3 we will see that a similar approach can be used with magnetism. The difference is that, whereas Gauss' law involves the flux, a measure of how much the field *spreads out*, the corresponding law for magnetism will measure how much the field *curls*.

Is it just dumb luck that the magnetic-field case came out so similar to the electric field case? Not at all. We've already seen that what one observer perceives as an electric field, another observer may perceive as a magnetic field. An observer flying along above a charged sheet will say that the charges are in motion, and will therefore say that it is both a sheet of current and a sheet of charge. Instead of a pure electric field, this observer will experience a combination of an electric field and a magnetic one. (We could also construct an example like the one in figure c on page 676, in which the field was purely magnetic.)

11.2.2 Energy in the magnetic field

In section 10.4, I've already argued that the energy density of the magnetic field must be proportional to $|\mathbf{B}|^2$, which we can write as B^2 for convenience. To pin down the constant of proportionality, we now need to do something like the argument on page 606: find one example where we can calculate the mechanical work done by the magnetic field, and equate that to the amount of energy lost by the field itself. The easiest example is two parallel sheets of charge, with their currents in opposite directions. Homework problem 53 is such a calculation, which gives the result

$$dU_m = \frac{c^2}{8\pi k} B^2 dv.$$

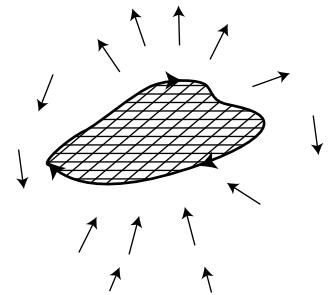
11.2.3 Superposition of dipoles

To understand this subsection, you'll have to have studied section 4.2.4, on iterated integrals.

The distant field of a dipole, in its midplane

Most current distributions cannot be broken down into long, straight wires, and subsection 11.2.1 has exhausted most of the interesting cases we can handle in this way. A much more useful building block is a square current loop. We have already seen how the dipole moment of an irregular current loop can be found by breaking the loop down into square dipoles (figure 1 on page 682), because the currents in adjoining squares cancel out on their shared edges. Likewise, as shown in figure f, if we could find the magnetic field of a square dipole, then we could find the field of any planar loop of current by adding the contributions to the field from all the squares.

The field of a square-loop dipole is very complicated close up, but luckily for us, we only need to know the current at distances that are large compared to the size of the loop, because we're free to make the squares on our grid as small as we like. The *distant* field of a square dipole turns out to be simple, and is no different from the distant field of any other dipole with the same dipole moment. We can also save ourselves some work if we only worry about finding the field of the dipole in its own plane, i.e., the plane perpendicular to its dipole moment. By symmetry, the field in this plane cannot have any component in the radial direction (inward toward the dipole, or outward away from it); it is perpendicular to the plane, and in the opposite direction compared to the dipole vector. (The field *inside* the loop is in the same direction as the dipole vector, but we're interested in the distant field.) Letting the dipole vector be along the z axis, we find that the field in the $x - y$ plane is of the form $B_z = f(r)$, where $f(r)$ is some function that depends only on r , the distance from the dipole.



f / The field of any planar current loop can be found by breaking it down into square dipoles.

We can pin down the result even more without any math. We know that the magnetic field made by a current always contains a factor of k/c^2 , which is the coupling constant for magnetism. We also know that the field must be proportional to the dipole moment, $m = IA$. Fields are always directly proportional to currents, and the proportionality to area follows because dipoles add according to their area. For instance, a square dipole that is 2 micrometers by 2 micrometers in size can be cut up into four dipoles that are 1 micrometer on a side. This tells us that our result must be of the form $B_z = (k/c^2)(IA)g(r)$. Now if we multiply the quantity $(k/c^2)(IA)$ by the function $g(r)$, we have to get units of teslas, and this only works out if $g(r)$ has units of m^{-3} (homework problem 15), so our result must be of the form

$$B_z = \frac{\beta kIA}{c^2 r^3},$$

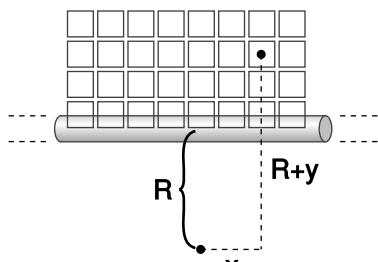
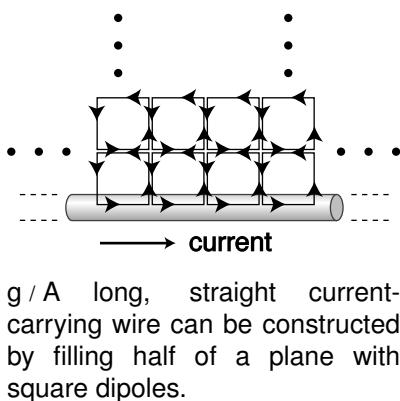
where β is a unitless constant. Thus our only task is to determine β , and we will have determined the field of the dipole (in the plane of its current, i.e., the midplane with respect to its dipole moment vector).

If we wanted to, we could simply build a dipole, measure its field, and determine β empirically. Better yet, we can get an exact result if we take a current loop whose field we know exactly, break it down into infinitesimally small squares, integrate to find the total field, set this result equal to the known expression for the field of the loop, and solve for β . There's just one problem here. We don't yet know an expression for the field of *any* current loop of *any* shape — all we know is the field of a long, straight wire. Are we out of luck? No, because, as shown in figure g, we can make a long, straight wire by putting together square dipoles! Any square dipole away from the edge has all four of its currents canceled by its neighbors. The only currents that don't cancel are the ones on the edge, so by superimposing all the square dipoles, we get a straight-line current.

This might seem strange. If the squares on the interior have all their currents canceled out by their neighbors, why do we even need them? Well, we need the squares on the edge in order to make the straight-line current. We need the second row of squares to cancel out the currents at the top of the first row of squares, and so on.

Integrating as shown in figure h, we have

$$B_z = \int_{y=0}^{\infty} \int_{x=-\infty}^{\infty} dB_z,$$



where $\mathrm{d}B_z$ is the contribution to the total magnetic field at our point of interest, which lies a distance R from the wire.

$$\begin{aligned} B_z &= \int_{y=0}^{\infty} \int_{x=-\infty}^{\infty} \frac{\beta k I \mathrm{d}A}{c^2 r^3} \\ &= \frac{\beta k I}{c^2} \int_{y=0}^{\infty} \int_{x=-\infty}^{\infty} \frac{1}{[x^2 + (R+y)^2]^{3/2}} \mathrm{d}x \mathrm{d}y \\ &= \frac{\beta k I}{c^2 R^3} \int_{y=0}^{\infty} \int_{x=-\infty}^{\infty} \left[\left(\frac{x}{R} \right)^2 + \left(1 + \frac{y}{R} \right)^2 \right]^{-3/2} \mathrm{d}x \mathrm{d}y \end{aligned}$$

This can be simplified with the substitutions $x = Ru$, $y = Rv$, and $\mathrm{d}x \mathrm{d}y = R^2 \mathrm{d}u \mathrm{d}v$:

$$B_z = \frac{\beta k I}{c^2 R} \int_{v=0}^{\infty} \int_{u=-\infty}^{\infty} \frac{1}{[u^2 + (1+v)^2]^{3/2}} \mathrm{d}u \mathrm{d}v$$

The u integral is of the form $\int_{-\infty}^{\infty} (u^2 + b)^{-3/2} \mathrm{d}u = 2/b^2$, so

$$B_z = \frac{\beta k I}{c^2 R} \int_{v=0}^{\infty} \frac{1}{(1+v)^2} \mathrm{d}v,$$

and the remaining v integral is equals 2, so

$$B_z = \frac{2\beta k I}{c^2 R}.$$

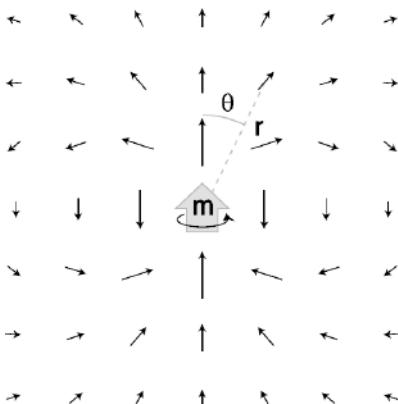
This is the field of a wire, which we already know equals $2kI/c^2R$, so we have $\beta=1$. Remember, the point of this whole calculation was not to find the field of a *wire*, which we already knew, but to find the unitless constant β in the expression for the field of a *dipole*. The distant field of a dipole, in its midplane, is therefore $B_z = \beta k I A / c^2 r^3 = k I A / c^2 r^3$, or, in terms of the dipole moment,

$$B_z = \frac{km}{c^2 r^3}.$$

The distant field of a dipole, out of its midplane

What about the field of a magnetic dipole outside of the dipole's midplane? Let's compare with an electric dipole. An electric dipole, unlike a magnetic one, can be built out of two opposite monopoles, i.e., charges, separated by a certain distance, and it is then straightforward to show by vector addition that the field of an electric dipole is

$$\begin{aligned} E_z &= kD (3 \cos^2 \theta - 1) r^{-3} \\ E_R &= kD (3 \sin \theta \cos \theta) r^{-3}, \end{aligned}$$



i / The field of a dipole.

where r is the distance from the dipole to the point of interest, θ is the angle between the dipole vector and the line connecting the dipole to this point, and E_z and E_R are, respectively, the components of the field parallel to and perpendicular to the dipole vector.

In the midplane, θ equals $\pi/2$, which produces $E_z = -kDr^{-3}$ and $E_R = 0$. This is the same as the field of a *magnetic* dipole in its midplane, except that the electric coupling constant k replaces the magnetic version k/c^2 , and the electric dipole moment D is substituted for the magnetic dipole moment m . It is therefore reasonable to conjecture that by using the same presto-change-o recipe we can find the field of a magnetic dipole outside its midplane:

$$B_z = \frac{km}{c^2} (3 \cos^2 \theta - 1) r^{-3}$$

$$B_R = \frac{km}{c^2} (3 \sin \theta \cos \theta) r^{-3}.$$

This turns out to be correct.⁵

Concentric, counterrotating currents

example 9

▷ Two concentric circular current loops, with radii a and b , carry the same amount of current I , but in opposite directions. What is the field at the center?

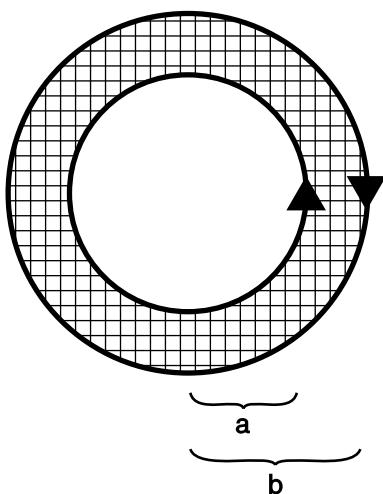
▷ We can produce these currents by tiling the region between the circles with square current loops, whose currents all cancel each other except at the inner and outer edges. The flavor of the calculation is the same as the one in which we made a line of current by filling a half-plane with square loops. The main difference is that this geometry has a different symmetry, so it will make more sense to use polar coordinates instead of x and y . The field at the center is

$$B_z = \int \frac{kI}{c^2 r^3} dA$$

$$= \int_{r=a}^b \frac{kI}{c^2 r^3} \cdot 2\pi r dr$$

$$= \frac{2\pi kI}{c^2} \left(\frac{1}{a} - \frac{1}{b} \right).$$

The positive sign indicates that the field is out of the page.



j / Example 9.

⁵If you've taken a course in differential equations, this won't seem like a very surprising assertion. The differential form of Gauss' law is a differential equation, and by giving the value of the field in the midplane, we've specified a boundary condition for the differential equation. Normally if you specify the boundary conditions, then there is a unique solution to the differential equation. In this particular case, it turns out that to ensure uniqueness, we also need to demand that the solution satisfy the differential form of Ampère's law, which is discussed in section 11.4.

Field at the center of a circular loop

example 10

▷ What is the magnetic field at the center of a circular current loop of radius a , which carries a current I ?

▷ This is like example 9, but with the outer loop being very large, and therefore too distant to make a significant field at the center. Taking the limit of that result as b approaches infinity, we have

$$B_z = \frac{2\pi k I}{c^2 a}$$

Comparing the results of examples 9 and 10, we see that the directions of the fields are both out of the page. In example 9, the outer loop has a current in the opposite direction, so it contributes a field that is into the page. This, however, is weaker than the field due to the inner loop, which dominates because it is less distant.

11.2.4 The g factor (optional)

In section 11.2.3 we exploited a particular trick for superimposing dipoles consisting of small square current loops. Let's now turn to a somewhat different way of superimposing dipoles. The idea is that matter is made out of atoms, which may act like little magnetic dipoles, but atoms are themselves made out of subatomic particles such as electrons, neutrons and protons — and there is no obvious way that we can ever know whether we have taken this process of reductionism (p. 18) to its conclusion. We can, however, look for clues in the electrical and mechanical properties of matter. Suppose that a particle of charge q and mass m is whizzing around and around some closed path. We don't even care whether the trajectory is a square or a circle, an orbit or a random wiggle. But let's say for convenience that it's a planar shape. The magnetic dipole moment (averaged over time) is $\mathbf{m} = IA$. But the angular momentum of a unit mass can also be interpreted (sec. 4.1.2, p. 256) as twice the area it sweeps out per unit time. Aside from the factor of two, which is just a historical glitch in the definitions, this mathematical analogy is exact: mass is to charge as angular momentum \mathbf{L} is to magnetic dipole moment \mathbf{m} . Therefore we have the identity

$$\frac{q}{m} \cdot \frac{|\mathbf{L}|}{|\mathbf{m}|} = 2$$

(where \mathbf{m} is the dipole moment, while m is the mass). The left-hand side is called the g factor. We expect $g = 2$ for a single orbiting particle.

Now suppose that we have a collection of particles with identical values of q/m (or a continuous distribution of charge and mass in which the ratio of the charge and mass densities is constant). Then vector addition of the \mathbf{L} and \mathbf{m} values gives the same $g = 2$ for the system as a whole. On the other hand, if the different members of

the system do *not* all have the same q/m , then the g of the system as a whole need not be 2. For example, a collection of positive and negative charges could easily have zero net charge but $\mathbf{m} \neq 0$, giving $g = 0$.

Particles such as the electron, the neutron, and the proton may be pointlike, or they may be composites of other particles. The electron and proton, which are charged, have the expected g factors of exactly 2 when we measure the \mathbf{L} and \mathbf{m} that they have due to their motion through space. But we also find that electrons, neutrons, and protons all come equipped with a built-in angular momentum, present even when they are at rest. This intrinsic angular momentum, called spin, is fixed in magnitude but can vary in direction, like that of a gyroscope. Thus if we measure the \mathbf{L} and \mathbf{m} of these particles when they are *at rest*, they have fixed g factors, which are as follows:

electron	2.002319304361
neutron	0
proton	5.58569471

The electron's intrinsic g factor is extremely close to 2, and if we ignore the small discrepancy for now, we are led to imagine that the electron is either a pointlike particle or a composite of smaller particles, each of which has the same charge-to-mass ratio. The neutron does have a nonvanishing dipole moment, so its zero g factor suggests that it is a composite of other particles whose charges cancel. The proton's g factor is quite different from 2, so we infer that it, too, is composite. The current theory is that protons and neutrons are clusters of particles called quarks. Quarks come in different types, and the different types have different values of q/m .

It is remarkable that we can infer these facts about the internal structures of neutrons and protons without having to do any experiments that directly probe their interior structure. We don't need a super-powerful microscope, nor do we need a particle accelerator that can supply enough energy to shake up their internal structure, like shaking a gift-wrapped box to tell what's inside. Merely by measuring the external, aggregate properties of the "box," we can get clues about the structure inside. This is closely analogous to the Tolman-Stewart experiment (example 11, p. 594), in which the subatomic structure of metals was probed by measuring inertial effects in an electric circuit. A more famous and important experiment using these ideas, by Stern and Gerlach, is described in sec. 14.1, p. 959.

11.2.5 The Biot-Savart law (optional)

In section 11.2.3 we developed a method for finding the field due to a given current distribution by tiling a plane with square dipoles. This method has several disadvantages:

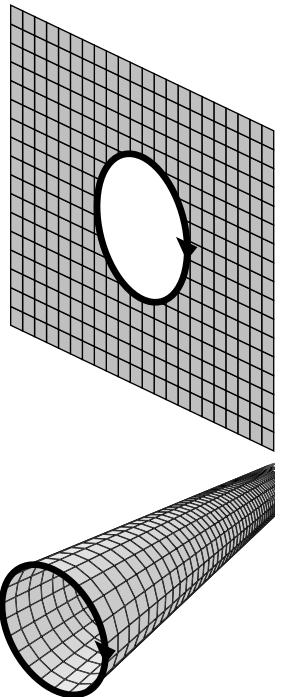
- The currents all have to lie in a single plane, and the point at which we're computing the field must be in that plane as well.
- We need to do integral over an area, which means one integral inside another, e.g., $\int \int \dots dx dy$. That can get messy.
- It's physically bizarre to have to construct square dipoles in places where there really aren't any currents.

Figure k shows the first step in eliminating these defects: instead of spreading our dipoles out in a plane, we bring them out along an axis. As shown in figure l, this eliminates the restriction to currents that lie in a plane. Now we have to use the general equations for a dipole field from page 696, rather than the simpler expression for the field in the midplane of a dipole. This increase in complication is more than compensated for by a fortunate feature of the new geometry, which is that the infinite tube can be broken down into strips, and we can find the field of such a strip for once and for all. This means that we no longer have to do one integral inside another. The derivation of the most general case is a little messy, so I'll just present the case shown in figure m, where the point of interest is assumed to lie in the $y - z$ plane. Intuitively, what we're really finding is the field of the short piece of length $d\ell$ on the end of the U; the two long parallel segments are going to be canceled out by their neighbors when we assemble all the strips to make the tube. We expect that the field of this end-piece will form a pattern that circulates around the y axis, so at the point of interest, it's really the x component of the field that we want to compute:

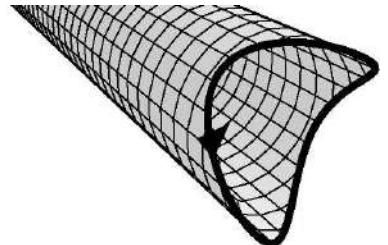
$$\begin{aligned} dB_x &= \int dB_R \cos \alpha \\ &= \int \frac{kI d\ell dx}{c^2 s^3} (3 \sin \theta \cos \theta \cos \alpha) \\ &= \frac{3kI d\ell}{c^2} \int_0^\infty \frac{1}{s^3} \left(\frac{xz}{s^2} \right) dx \\ &= \frac{3kI z d\ell}{c^2} \int_0^\infty \frac{x}{(x^2 + r^2)^{5/2}} dx \\ &= \frac{kI d\ell z}{c^2 r^3} \\ &= \frac{kI d\ell \sin \phi}{c^2 r^2} \end{aligned}$$

In the more general case, l, the current loop is not planar, the point of interest is not in the end-planes of the U's, and the U shapes have their ends staggered, so the end-piece $d\ell$ is not the only part of each U whose current is not canceled. Without going into the gory details, the correct general result is as follows:

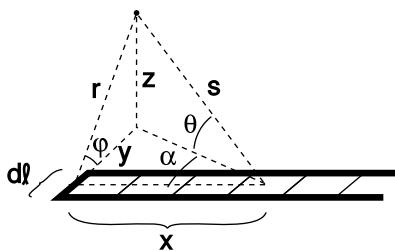
$$dB = \frac{kI d\ell \times \mathbf{r}}{c^2 r^3},$$



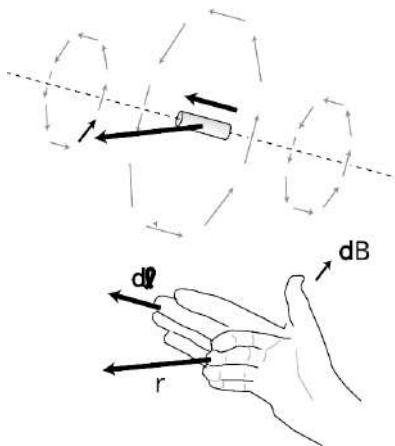
k / Two ways of making a current loop out of square dipoles.



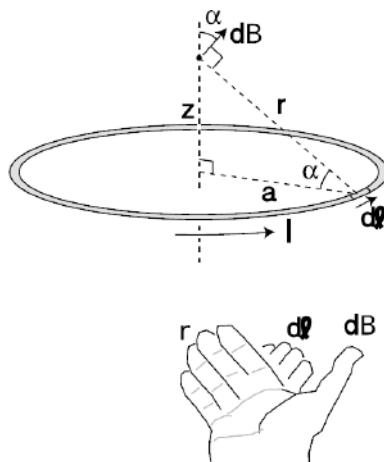
l / The new method can handle non-planar currents.



m / The field of an infinite U.



n / The geometry of the Biot-Savart law. The small arrows show the result of the Biot-Savart law at various positions relative to the current segment $d\ell$. The Biot-Savart law involves a cross product, and the right-hand rule for this cross product is demonstrated for one case.



o / Example 12.

which is known as the Biot-Savart law. (It rhymes with “leo bazaar.” Both t’s are silent.) The distances $d\ell$ and r are now defined as vectors, $d\ell$ and \mathbf{r} , which point, respectively, in the direction of the current in the end-piece and the direction from the end-piece to the point of interest. The new equation looks different, but it is consistent with the old one. The vector cross product $d\ell \times \mathbf{r}$ has a magnitude $r d\ell \sin \phi$, which cancels one of r ’s in the denominator and makes the $d\ell \times \mathbf{r}/r^3$ into a vector with magnitude $d\ell \sin \phi/r^2$.

The field at the center of a circular loop

example 11

Previously we had to do quite a bit of work (examples 9 and 10), to calculate the field at the center of a circular loop of current of radius a . It’s much easier now. Dividing the loop into many short segments, each $d\ell$ is perpendicular to the \mathbf{r} vector that goes from it to the center of the circle, and every \mathbf{r} vector has magnitude a . Therefore every cross product $d\ell \times \mathbf{r}$ has the same magnitude, $a d\ell$, as well as the same direction along the axis perpendicular to the loop. The field is

$$\begin{aligned} B &= \int \frac{k I a d\ell}{c^2 a^3} \\ &= \frac{k I}{c^2 a^2} \int d\ell \\ &= \frac{k I}{c^2 a^2} (2\pi a) \\ &= \frac{2\pi k I}{c^2 a} \end{aligned}$$

Out-of-the-plane field of a circular loop

example 12

▷ What is the magnetic field of a circular loop of current at a point on the axis perpendicular to the loop, lying a distance z from the loop’s center?

▷ Again, let’s write a for the loop’s radius. The \mathbf{r} vector now has magnitude $\sqrt{a^2 + z^2}$, but it is still perpendicular to the $d\ell$ vector. By symmetry, the only nonvanishing component of the field is along the z axis,

$$\begin{aligned} B_z &= \int |d\mathbf{B}| \cos \alpha \\ &= \int \frac{k I r d\ell}{c^2 r^3} \frac{a}{r} \\ &= \frac{k I a}{c^2 r^3} \int d\ell \\ &= \frac{2\pi k I a^2}{c^2 (a^2 + z^2)^{3/2}}. \end{aligned}$$

Is it the field of a particle?

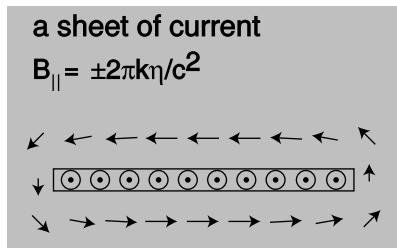
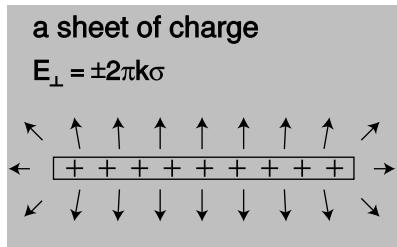
We have a simple equation, based on Coulomb's law, for the electric field surrounding a charged particle. Looking at figure n, we can imagine that if the current segment $d\ell$ was very short, then it might only contain one electron. It's tempting, then, to interpret the Biot-Savart law as a similar equation for the magnetic field surrounding a moving charged particle. Tempting but wrong! Suppose you stand at a certain point in space and watch a charged particle move by. It has an electric field, and since it's moving, you will also detect a magnetic field on top of that. Both of these fields change over time, however. Not only do they change their magnitudes and directions due to your changing geometric relationship to the particle, but they are also time-delayed, because disturbances in the electromagnetic field travel at the speed of light, which is finite. The fields you detect are the ones corresponding to where the particle used to be, not where it is now. Coulomb's law and the Biot-Savart law are both false in this situation, since neither equation includes time as a variable. It's valid to think of Coulomb's law as the equation for the field of a stationary charged particle, but not a moving one. The Biot-Savart law fails completely as a description of the field of a charged particle, since stationary particles don't make magnetic fields, and the Biot-Savart law fails in the case where the particle is moving.

If you look back at the long chain of reasoning that led to the Biot-Savart law, it all started from the relativistic arguments at the beginning of this chapter, where we assumed a steady current in an infinitely long wire. Everything that came later was built on this foundation, so all our reasoning depends on the assumption that the currents are steady. In a steady current, any charge that moves away from a certain spot is replaced by more charge coming up behind it, so even though the charges are all moving, the electric and magnetic fields they produce are constant. Problems of this type are called electrostatics and magnetostatics problems, and it is only for these problems that Coulomb's law and the Biot-Savart law are valid.

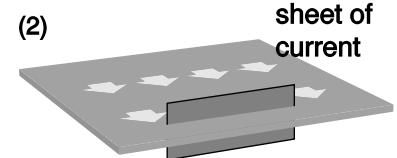
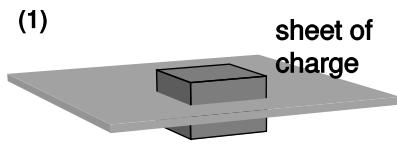
You might think that we could patch up Coulomb's law and the Biot-Savart law by inserting the appropriate time delays. However, we've already seen a clear example of a phenomenon that wouldn't be fixed by this patch: on page 622, we found that a changing magnetic field creates an electric field. Induction effects like these also lead to the existence of light, which is a wave disturbance in the electric and magnetic fields. We could try to apply another band-aid fix to Coulomb's law and the Biot-Savart law to make them deal with induction, but it won't work.

So what *are* the fundamental equations that describe how sources give rise to electromagnetic fields? We've already encountered two of them: Gauss' law for electricity and Gauss' law for magnetism.

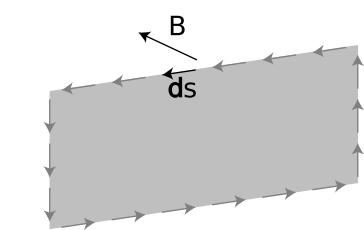
Experiments show that these are valid in all situations, not just static ones. But Gauss' law for magnetism merely says that the magnetic flux through a closed surface is zero. It doesn't tell us how to make magnetic fields using currents. It only tells us that we *can't* make them using magnetic monopoles. The following section develops a new equation, called Ampère's law, which is equivalent to the Biot-Savart law for magnetostatics, but which, unlike the Biot-Savart law, can easily be extended to nonstatic situations.



a / The electric field of a sheet of charge, and the magnetic field of a sheet of current.



b / A Gaussian surface and an Ampèrian surface.



c / The definition of the circulation, Γ .

11.3 Magnetic fields by Ampère's law

11.3.1 Ampère's law

As discussed at the end of subsection 11.2.5, our goal now is to find an equation for magnetism that, unlike the Biot-Savart law, will not end up being a dead end when we try to extend it to nonstatic situations.⁶ Experiments show that Gauss' law is valid in both static and nonstatic situations, so it would be reasonable to look for an approach to magnetism that is similar to the way Gauss' law deals with electricity.

How can we do this? Figure a, reproduced from page 692, is our roadmap. Electric fields spread out from charges. Magnetic fields curl around currents. In figure b/1, we define a Gaussian surface, and we define the flux in terms of the electric field pointing out through this surface. In the magnetic case, b/2, we define a surface, called an Ampèrian surface, and we define a quantity called the circulation, Γ (uppercase Greek gamma), in terms of the magnetic field that points along the edge of the Ampèrian surface, c. We break the edge into tiny parts \mathbf{s}_j , and for each of these parts, we define a contribution to the circulation using the dot product of $d\mathbf{s}$ with the magnetic field:

$$\Gamma = \sum \mathbf{s}_j \cdot \mathbf{B}_j$$

The circulation is a measure of how curly the field is. Like a Gaussian surface, an Ampèrian surface is purely a mathematical construction. It is not a physical object.

In figure b/2, the field is perpendicular to the edges on the ends, but parallel to the top and bottom edges. A dot product is zero when the vectors are perpendicular, so only the top and bottom edges contribute to Γ . Let these edges have length s . Since the field is constant along both of these edges, we don't actually have to break them into tiny parts; we can just have \mathbf{s}_1 on the top edge, pointing to the left, and \mathbf{s}_2 on the bottom edge, pointing to the right. The vector \mathbf{s}_1 is in the same direction as the field \mathbf{B}_1 , and \mathbf{s}_2 is in the same direction as \mathbf{B}_2 , so the dot products are simply equal to

⁶If you didn't read this optional subsection, don't worry, because the point is that we need to try a whole new approach anyway.

the products of the vectors' magnitudes. The resulting circulation is

$$\begin{aligned}\Gamma &= |\mathbf{s}_1||\mathbf{B}_1| + |\mathbf{s}_2||\mathbf{B}_2| \\ &= \frac{2\pi k\eta s}{c^2} + \frac{2\pi k\eta s}{c^2} \\ &= \frac{4\pi k\eta s}{c^2}.\end{aligned}$$

But ηs is (current/length)(length), i.e., it is the amount of current that pierces the Ampèrean surface. We'll call this current I_{through} . We have found one specific example of the general law of nature known as Ampère's law:

$$\Gamma = \frac{4\pi k}{c^2} I_{\text{through}}$$

Positive and negative signs

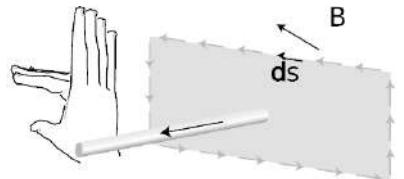
Figures d/1 and d/2 show what happens to the circulation when we reverse the direction of the current I_{through} . Reversing the current causes the magnetic field to reverse itself as well. The dot products occurring in the circulation are all negative in d/2, so the total circulation is now negative. To preserve Ampère's law, we need to define the current in d/2 as a negative number. In general, determine these plus and minus signs using the right-hand rule shown in the figure. As the fingers of your hand sweep around in the direction of the \mathbf{s} vectors, your thumb defines the direction of current which is positive. Choosing the direction of the thumb is like choosing which way to insert an ammeter in a circuit: on a digital meter, reversing the connections gives readings which are opposite in sign.

A solenoid

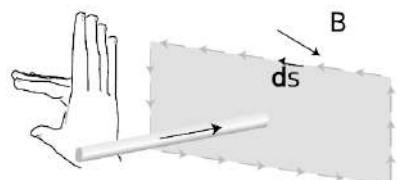
example 13

▷ What is the field inside a long, straight solenoid of length ℓ and radius a , and having N loops of wire evenly wound along it, which carry a current I ?

▷ This is an interesting example, because it allows us to get a very good approximation to the field, but without some experimental input it wouldn't be obvious what approximation to use. Figure e/1 shows what we'd observe by measuring the field of a real solenoid. The field is nearly constant inside the tube, as long as we stay far away from the mouths. The field outside is much weaker. For the sake of an approximate calculation, we can idealize this field as shown in figure e/2. Of the edges of the Ampèrean surface shown in e/3, only AB contributes to the flux — there is zero field along CD, and the field is perpendicular to edges BC

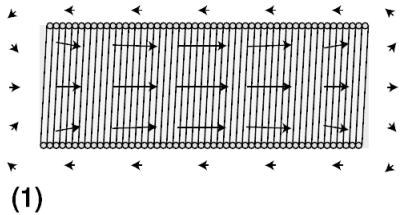


(1) $\Gamma > 0, I_{\text{through}} > 0$

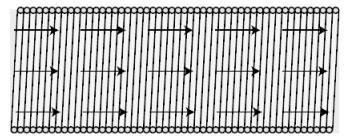


(2) $\Gamma < 0, I_{\text{through}} < 0$

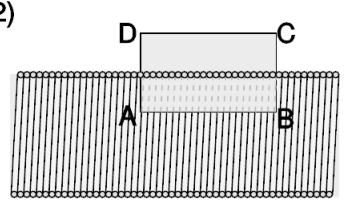
d / Positive and negative signs in Ampère's law.



(1)



(2)



(3)

e / Example 13: a cutaway view of a solenoid.

and DA. Ampère's law gives

$$\begin{aligned}\Gamma &= \frac{4\pi k}{c^2} I_{\text{through}} \\ (B)(\text{length of AB}) &= \frac{4\pi k}{c^2} (\eta)(\text{length of AB}) \\ B &= \frac{4\pi k\eta}{c^2} \\ &= \frac{4\pi kNI}{c^2\ell}\end{aligned}$$

self-check D

What direction is the current in figure e?

▷ Answer, p. 1064

self-check E

Based on how ℓ entered into the derivation in example 13, how should it be interpreted? Is it the total length of the wire? ▷ Answer, p. 1065

self-check F

Surprisingly, we never needed to know the radius of the solenoid in example 13. Why is it physically plausible that the answer would be independent of the radius? ▷ Answer, p. 1065

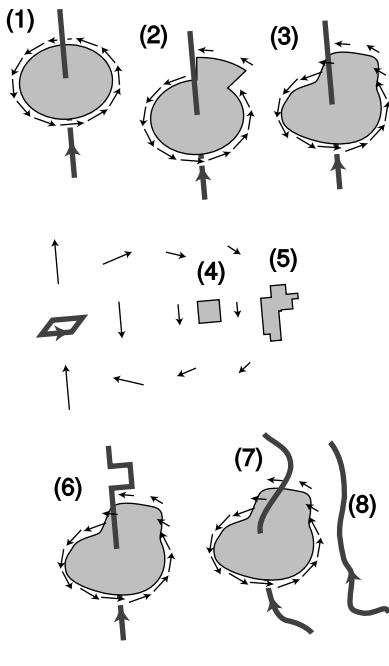
Example 13 shows how much easier it can sometimes be to calculate a field using Ampère's law rather than the approaches developed previously in this chapter. However, if we hadn't already known something about the field, we wouldn't have been able to get started. In situations that lack symmetry, Ampère's law may make things harder, not easier. Anyhow, we will have no choice in nonstatic cases, where Ampère's law is true, and static equations like the Biot-Savart law are false.

11.3.2 A quick and dirty proof

Here's an informal sketch for a proof of Ampère's law, with no pretensions to rigor. Even if you don't care much for proofs, it would be a good idea to read it, because it will help to build your ability to visualize how Ampère's law works.

First we establish by a direct computation (homework problem 26) that Ampère's law holds for the geometry shown in figure f/1, a circular Ampèrian surface with a wire passing perpendicularly through its center. If we then alter the surface as in figure f/2, Ampère's law still works, because the straight segments, being perpendicular to the field, don't contribute to the circulation, and the new arc makes the same contribution to the circulation as the old one it replaced, because the weaker field is compensated for by the greater length of the arc. It is clear that by a series of such modifications, we could mold the surface into any shape, f/3.

Next we prove Ampère's law in the case shown in figure f/4: a small, square Ampèrian surface subject to the field of a distant



f / A proof of Ampère's law.

square dipole. This part of the proof can be most easily accomplished by the methods of section 11.4. It should, for example, be plausible in the case illustrated here. The field on the left edge is stronger than the field on the right, so the overall contribution of these two edges to the circulation is slightly counterclockwise. However, the field is not quite perpendicular to the top and bottoms edges, so they both make small clockwise contributions. The clockwise and counterclockwise parts of the circulation end up canceling each other out. Once Ampère's law is established for a square surface like $f/4$, it follows that it is true for an irregular surface like $f/5$, since we can build such a shape out of squares, and the circulations are additive when we paste the surfaces together this way.

By pasting a square dipole onto the wire, $f/6$, like a flag attached to a flagpole, we can cancel out a segment of the wire's current and create a detour. Ampère's law is still true because, as shown in the last step, the square dipole makes zero contribution to the circulation. We can make as many detours as we like in this manner, thereby morphing the wire into an arbitrary shape like $f/7$.

What about a wire like $f/8$? It doesn't pierce the Ampèreian surface, so it doesn't add anything to I_{through} , and we need to show that it likewise doesn't change the circulation. This wire, however, could be built by tiling the half-plane on its right with square dipoles, and we've already established that the field of a distant dipole doesn't contribute to the circulation. (Note that we couldn't have done this with a wire like $f/7$, because some of the dipoles would have been right on top of the Ampèreian surface.)

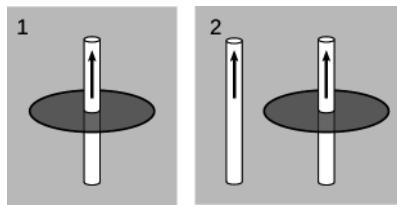
If Ampère's law holds for cases like $f/7$ and $f/8$, then it holds for any complex bundle of wires, including some that pass through the Ampèreian surface and some that don't. But we can build just about any static current distribution we like using such a bundle of wires, so it follows that Ampère's law is valid for any static current distribution.

11.3.3 Maxwell's equations for static fields

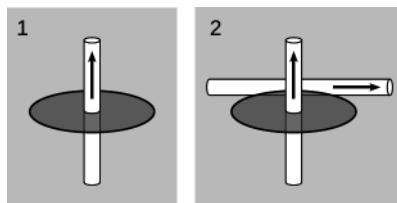
Static electric fields *don't* curl the way magnetic fields do, so we can state a version of Ampère's law for electric fields, which is that the circulation of the electric field is zero. Summarizing what we know so far about static fields, we have

$$\begin{aligned}\Phi_E &= 4\pi k q_{in} \\ \Phi_B &= 0 \\ \Gamma_E &= 0 \\ \Gamma_B &= \frac{4\pi k}{c^2} I_{\text{through}}.\end{aligned}$$

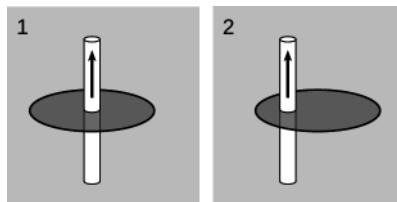
This set of equations is the static case of the more general relations known as Maxwell's equations. On the left side of each equation, we



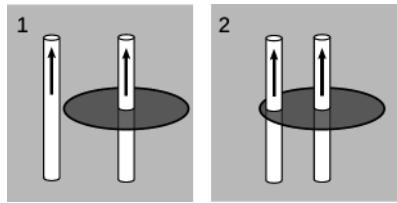
g / Discussion question A.



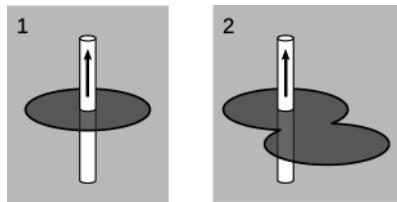
h / Discussion question B.



i / Discussion question C.



j / Discussion question D.



k / Discussion question E.

have information about a field. On the right is information about the field's sources.

It is vitally important to realize that these equations are only true for statics. They are incorrect if the distribution of charges or currents is changing over time. For example, we saw on page 622 that the changing magnetic field in an inductor gives rise to an electric field. Such an effect is completely inconsistent with the static version of Maxwell's equations; the equations don't even refer to time, so if the magnetic field is changing over time, they will not do anything special. The extension of Maxwell's equations to nonstatic fields is discussed in section 11.6.

Discussion Questions

A Figure g/1 shows a wire with a circular Ampèrean surface drawn around its waist; in this situation, Ampère's law can be verified easily based on the equation for the field of a wire. In panel 2, a second wire has been added. Explain why it's plausible that Ampère's law still holds.

B Figure h is like figure g, but now the second wire is perpendicular to the first, and lies in the plane of, and outside of, the Ampèrean surface. Carry out a similar analysis.

C This discussion question is similar to questions A and B, but now the Ampèrean surface has been moved off center.

D The left-hand wire has been nudged over a little. Analyze as before.

E You know what to do.

11.4 Ampère's law in differential form (optional)

11.4.1 The curl operator

The differential form of Gauss' law is more physically satisfying than the integral form, because it relates the charges that are present at some point to the properties of the electric field *at the same point*. Likewise, it would be more attractive to have a differential version of Ampère's law that would relate the currents to the magnetic field at a single point. Intuitively, the divergence was based on the idea of the div-meter, a/1. The corresponding device for measuring the curliness of a field is the curl-meter, a/2. If the field is curly, then the torques on the charges will not cancel out, and the wheel will twist against the resistance of the spring. If your intuition tells you that the curlmeter will never do anything at all, then your intuition is doing a beautiful job on static fields; for nonstatic fields, however, it is perfectly possible to get a curly electric field.

Gauss' law in differential form relates $\text{div } \mathbf{E}$, a scalar, to the charge density, another scalar. Ampère's law, however, deals with directions in space: if we reverse the directions of the currents, it makes a difference. We therefore expect that the differential form of Ampère's law will have vectors on both sides of the equal sign, and we should be thinking of the curl-meter's result as a vector. First we find the orientation of the curl-meter that gives the strongest torque, and then we define the direction of the curl vector using the right-hand rule shown in figure a/3.

To convert the div-meter concept to a mathematical definition, we found the infinitesimal flux, $d\Phi$ through a tiny cubical Gaussian surface containing a volume dv . By analogy, we imagine a tiny square Ampèrian surface with infinitesimal area $d\mathbf{A}$. We assume this surface has been oriented in order to get the maximum circulation. The area vector $d\mathbf{A}$ will then be in the same direction as the one defined in figure a/3. Ampère's law is

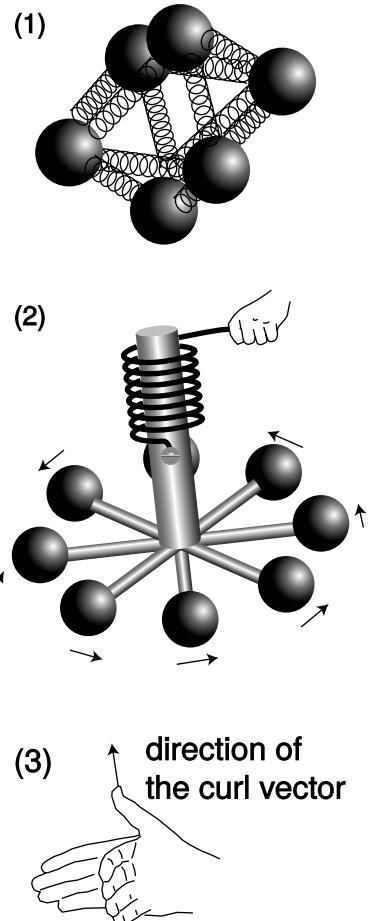
$$d\Gamma = \frac{4\pi k}{c^2} dI_{\text{through}}.$$

We define a current density per unit area, \mathbf{j} , which is a vector pointing in the direction of the current and having magnitude $\mathbf{j} = dI/|d\mathbf{A}|$. In terms of this quantity, we have

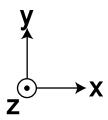
$$\begin{aligned} d\Gamma &= \frac{4\pi k}{c^2} j |\mathbf{j}| |d\mathbf{A}| \\ \frac{d\Gamma}{|d\mathbf{A}|} &= \frac{4\pi k}{c^2} |\mathbf{j}| \end{aligned}$$

With this motivation, we define the magnitude of the curl as

$$|\text{curl } \mathbf{B}| = \frac{d\Gamma}{|d\mathbf{A}|}.$$

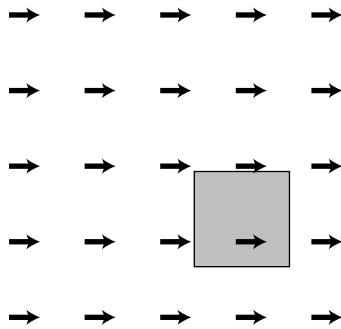


a / The div-meter, 1, and the curl-meter, 2 and 3.

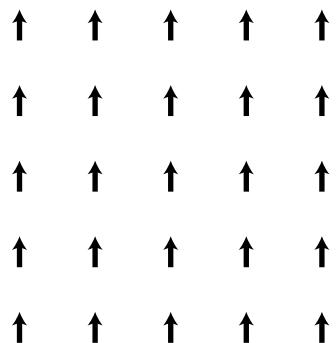


Note that the curl, just like a derivative, has a differential divided by another differential. In terms of this definition, we find Ampère's law in differential form:

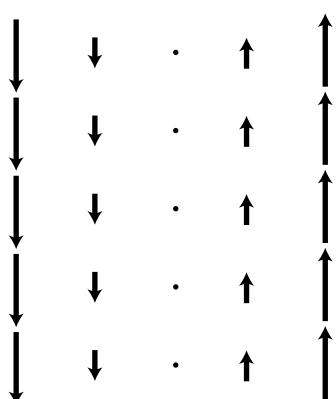
b / The coordinate system used in the following examples.



c / The field \hat{x} .



d / The field \hat{y} .



e / The field $\hat{x}\hat{y}$.

$$\text{curl } \mathbf{B} = \frac{4\pi k}{c^2} \mathbf{j}$$

The complete set of Maxwell's equations in differential form is collected on page 1029.

11.4.2 Properties of the curl operator

The curl is a derivative.

As an example, let's calculate the curl of the field \hat{x} shown in figure c. For our present purposes, it doesn't really matter whether this is an electric or a magnetic field; we're just getting our feet wet with the curl as a mathematical definition. Applying the definition of the curl directly, we construct an Amperian surface in the shape of an infinitesimally small square. Actually, since the field is uniform, it doesn't even matter very much whether we make the square finite or infinitesimal. The right and left edges don't contribute to the circulation, since the field is perpendicular to these edges. The top and bottom do contribute, but the top's contribution is clockwise, i.e., into the page according to the right-hand rule, while the bottom contributes an equal amount in the counterclockwise direction, which corresponds to an out-of-the-page contribution to the curl. They cancel, and the circulation is zero. We could also have determined this by imagining a curl-meter inserted in this field: the torques on it would have canceled out.

It makes sense that the curl of a constant field is zero, because the curl is a kind of derivative. The derivative of a constant is zero.

The curl is rotationally invariant.

Figure c looks just like figure c, but rotated by 90 degrees. Physically, we could be viewing the same field from a point of view that was rotated. Since the laws of physics are the same regardless of rotation, the curl must be zero here as well. In other words, the curl is rotationally invariant. If a certain field has a certain curl vector, then viewed from some other angle, we simply see the same field and the same curl vector, viewed from a different angle. A zero vector viewed from a different angle is still a zero vector.

As a less trivial example, let's compute the curl of the field $\mathbf{F} = x\hat{y}$ shown in figure e, at the point $(x = 0, y = 0)$. The circulation around a square of side s centered on the origin can be approximated

by evaluating the field at the midpoints of its sides,

$$\begin{array}{llll} x = s/2 & y = 0 & \mathbf{F} = (s/2)\hat{\mathbf{y}} & \mathbf{s}_1 \cdot \mathbf{F} = s^2/2 \\ x = 0 & y = s/2 & \mathbf{F} = 0 & \mathbf{s}_2 \cdot \mathbf{F} = 0 \\ x = -s/2 & y = 0 & \mathbf{F} = -(s/2)\hat{\mathbf{y}} & \mathbf{s}_3 \cdot \mathbf{F} = s^2/2 \\ x = 0 & y = -s/2 & \mathbf{F} = 0 & \mathbf{s}_4 \cdot \mathbf{F} = 0, \end{array}$$

which gives a circulation of s^2 , and a curl with a magnitude of $s^2/\text{area} = s^2/s^2 = 1$. By the right-hand rule, the curl points out of the page, i.e., along the positive z axis, so we have

$$\text{curl } xy\hat{\mathbf{y}} = \hat{\mathbf{z}}.$$

Now consider the field $-y\hat{\mathbf{x}}$, shown in figure f. This is the same as the previous field, but with your book rotated by 90 degrees about the z axis. Rotating the result of the first calculation, $\hat{\mathbf{z}}$, about the z axis doesn't change it, so the curl of this field is also $\hat{\mathbf{z}}$.

Scaling

When you're taking an ordinary derivative, you have the rule

$$\frac{d}{dx}[cf(x)] = c\frac{d}{dx}f(x).$$

In other words, multiplying a function by a constant results in a derivative that is multiplied by that constant. The curl is a kind of derivative operator, and the same is true for a curl.

Multiplying the field by -1 .

example 14

▷ What is the curl of the field $-x\hat{\mathbf{y}}$ at the origin?

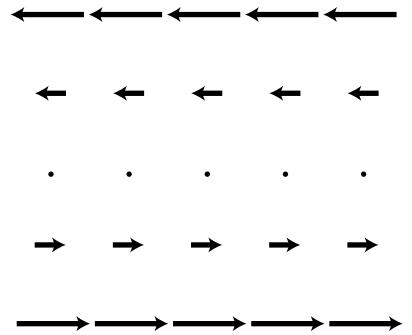
▷ Using the scaling property just discussed, we can make this into a curl that we've already calculated:

$$\begin{aligned} \text{curl } (-x\hat{\mathbf{y}}) &= -\text{curl } (x\hat{\mathbf{y}}) \\ &= -\hat{\mathbf{z}} \end{aligned}$$

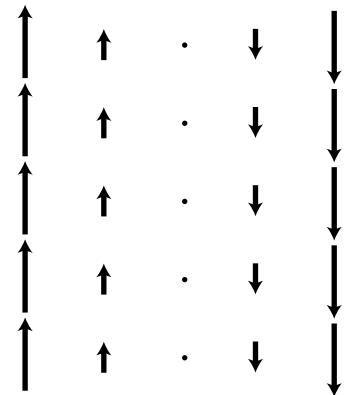
This is in agreement with the right-hand rule.

The curl is additive.

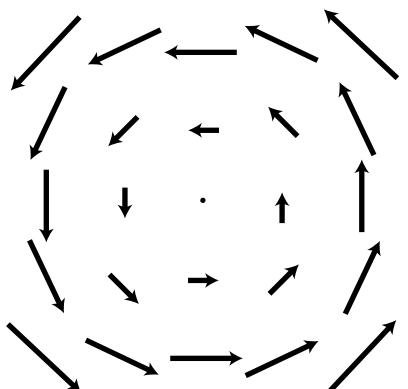
We have only calculated each field's curl at the origin, but each of these fields actually has the same curl everywhere. In example 14, for instance, it is obvious that the curl is constant along any vertical line. But even if we move along the x axis, there is still an imbalance between the torques on the left and right sides of the curl-meter. More formally, suppose we start from the origin and move to the left by one unit. We find ourselves in a region where the field



f / The field $-y\hat{\mathbf{x}}$.



g / Example 14.



h / Example 15.

is very much as it was before, except that all the field vectors have had one unit worth of $\hat{\mathbf{y}}$ added to them. But what do we get if we take the curl of $-x\hat{\mathbf{y}} + \hat{\mathbf{y}}$? The curl, like any god-fearing derivative operation, has the additive property

$$\operatorname{curl} (\mathbf{F} + \mathbf{G}) = \operatorname{curl} \mathbf{F} + \operatorname{curl} \mathbf{G},$$

so

$$\operatorname{curl} (-x\hat{\mathbf{y}} + \hat{\mathbf{y}}) = \operatorname{curl} (-x\hat{\mathbf{y}}) + \operatorname{curl} (\hat{\mathbf{y}}).$$

But the second term is zero, so we get the same result as at the origin.

A field that goes in a circle

example 15

▷ What is the curl of the field $x\hat{\mathbf{y}} - y\hat{\mathbf{x}}$?

▷ Using the linearity of the curl, and recognizing each of the terms as one whose curl we have already computed, we find that this field's curl is a constant $2\hat{\mathbf{z}}$. This agrees with the right-hand rule.

The field inside a long, straight wire

example 16

▷ What is the magnetic field *inside* a long, straight wire in which the current density is j ?

▷ Let the wire be along the z axis, so $\mathbf{j} = j\hat{\mathbf{z}}$. Ampère's law gives

$$\operatorname{curl} \mathbf{B} = \frac{4\pi k}{c^2} j\hat{\mathbf{z}}.$$

In other words, we need a magnetic field whose curl is a constant. We've encountered several fields with constant curls, but the only one that has the same symmetry as the cylindrical wire is $x\hat{\mathbf{y}} - y\hat{\mathbf{x}}$, so the answer must be this field or some constant multiplied by it,

$$\mathbf{B} = b(x\hat{\mathbf{y}} - y\hat{\mathbf{x}}).$$

The curl of this field is $2b\hat{\mathbf{z}}$, so

$$2b = \frac{4\pi k}{c^2} j,$$

and thus

$$\mathbf{B} = \frac{2\pi k}{c^2} j(x\hat{\mathbf{y}} - y\hat{\mathbf{x}}).$$

The curl in component form

Now consider the field

$$F_x = ax + by + c$$

$$F_y = dx + ey + f,$$

i.e.,

$$\mathbf{F} = ax\hat{\mathbf{x}} + by\hat{\mathbf{x}} + c\hat{\mathbf{x}} + dx\hat{\mathbf{y}} + ey\hat{\mathbf{y}} + f\hat{\mathbf{y}}.$$

The only terms whose curls we haven't yet explicitly computed are the a , e , and f terms, and their curls turn out to be zero (homework problem 50). Only the b and d terms have nonvanishing curls. The curl of this field is

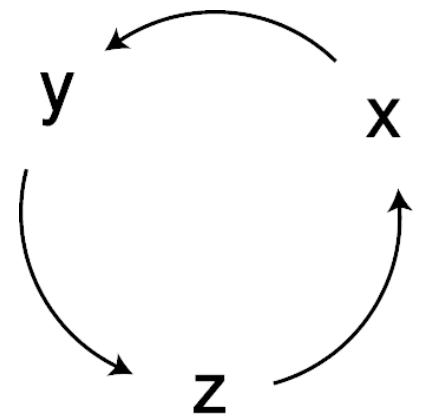
$$\begin{aligned}\operatorname{curl} \mathbf{F} &= \operatorname{curl}(by\hat{\mathbf{x}}) + \operatorname{curl}(dx\hat{\mathbf{y}}) \\ &= b \operatorname{curl}(y\hat{\mathbf{x}}) + d \operatorname{curl}(x\hat{\mathbf{y}}) \quad [\text{scaling}] \\ &= b(-\hat{\mathbf{z}}) + d(\hat{\mathbf{z}}) \quad [\text{found previously}] \\ &= (d - b)\hat{\mathbf{z}}.\end{aligned}$$

But *any* field in the $x - y$ plane can be approximated with this type of field, as long as we only need to get a good approximation within a small region. The infinitesimal Ampèrean surface occurring in the definition of the curl is tiny enough to fit in a pretty small region, so we can get away with this here. The d and b coefficients can then be associated with the partial derivatives $\partial F_y / \partial x$ and $\partial F_x / \partial y$. We therefore have

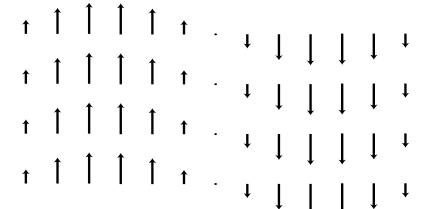
$$\operatorname{curl} \mathbf{F} = \left(\frac{\partial F_y}{\partial x} - \frac{\partial F_x}{\partial y} \right) \hat{\mathbf{z}}$$

for any field in the $x - y$ plane. In three dimensions, we just need to generate two more equations like this by doing a cyclic permutation of the variables x , y , and z :

$$\begin{aligned}(\operatorname{curl} \mathbf{F})_x &= \frac{\partial F_z}{\partial y} - \frac{\partial F_y}{\partial z} \\ (\operatorname{curl} \mathbf{F})_y &= \frac{\partial F_x}{\partial z} - \frac{\partial F_z}{\partial x} \\ (\operatorname{curl} \mathbf{F})_z &= \frac{\partial F_y}{\partial x} - \frac{\partial F_x}{\partial y}\end{aligned}$$



i / A cyclic permutation of x , y , and z .



j / Example 17.

A sine wave

example 17

▷ Find the curl of the following electric field

$$\mathbf{E} = (\sin x)\hat{\mathbf{y}},$$

and interpret the result.

▷ The only nonvanishing partial derivative occurring in this curl is

$$\begin{aligned}(\operatorname{curl} \mathbf{E})_z &= \frac{\partial E_y}{\partial x} \\ &= \cos x,\end{aligned}$$

so

$$\operatorname{curl} \mathbf{E} = \cos \hat{\mathbf{z}}$$

This is visually reasonable: the curl-meter would spin if we put its wheel in the plane of the page, with its axle poking out along the z axis. In some areas it would spin clockwise, in others counter-clockwise, and this makes sense, because the cosine is positive in some places and negative in others.

This is a perfectly reasonable field pattern: it's the electric field pattern of a light wave! But Ampère's law for electric fields says the curl of \mathbf{E} is supposed to be zero. What's going on? What's wrong is that we can't assume the *static* version of Ampère's law. All we've really proved is that this pattern is impossible as a static field: we can't have a light wave that stands still.

Figure k is a summary of the vector calculus presented in the optional sections of this book. The first column shows that one function is related to another by a kind of differentiation. The second column states the fundamental theorem of calculus, which says that if you integrate the derivative over the interior of a region, you get some information about the original function at the boundary of that region.

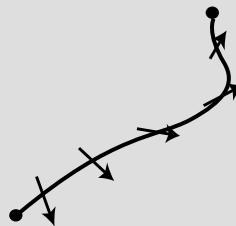
$$f = \frac{dg}{dx}$$

$$\Delta g = \int_{\text{interior}} dx$$



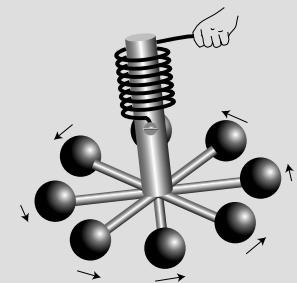
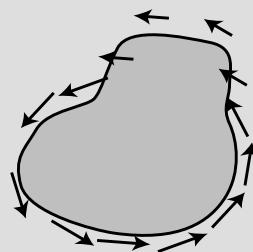
$$f = \nabla g$$

$$\Delta g = \int_{\text{interior}} ds$$



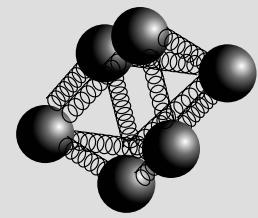
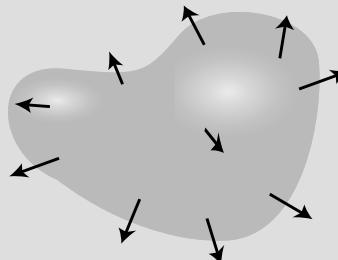
$$f = \operatorname{curl} g$$

$$\oint g = \int_{\text{interior}} f dA$$



$$f = \operatorname{div} g$$

$$\oint g = \int_{\text{interior}} f dv$$



k / A summary of the derivative, gradient, curl, and divergence.

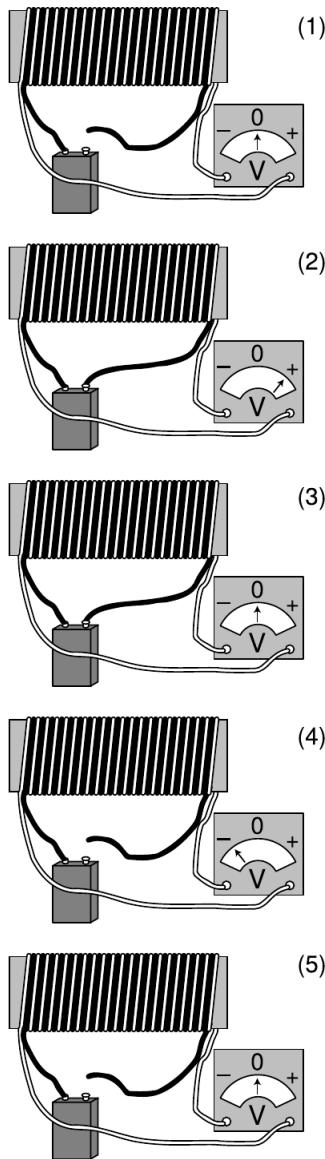
11.5 Induced electric fields

11.5.1 Faraday's experiment

Nature is simple, but the simplicity may not become evident until a hundred years after the discovery of some new piece of physics. We've already seen, on page 622, that the time-varying magnetic field in an inductor causes an electric field. This electric field is *not created by charges*. That argument, however, only seems clear with hindsight. The discovery of this phenomenon of induced electric fields — fields that are not due to charges — was a purely experimental accomplishment by Michael Faraday (1791-1867), the son of a blacksmith who had to struggle against the rigid class struc-



a / Faraday on a British banknote.



b / Faraday's experiment, simplified and shown with modern equipment.

ture of 19th century England. Faraday, working in 1831, had only a vague and general idea that electricity and magnetism were related to each other, based on Oersted's demonstration, a decade before, that magnetic fields were caused by electric currents.

Figure b is a simplified drawing of the following experiment, as described in Faraday's original paper: "Two hundred and three feet of copper wire . . . were passed round a large block of wood; [another] two hundred and three feet of similar wire were interposed as a spiral between the turns of the first, and metallic contact everywhere prevented by twine [insulation]. One of these [coils] was connected with a galvanometer [voltmeter], and the other with a battery. . . When the contact was made, there was a sudden and very slight effect at the galvanometer, and there was also a similar slight effect when the contact with the battery was broken. But whilst the . . . current was continuing to pass through the one [coil], no . . . effect . . . upon the other [coil] could be perceived, although the active power of the battery was proved to be great, by its heating the whole of its own coil [through ordinary resistive heating] . . ."

From Faraday's notes and publications, it appears that the situation in figure b/3 was a surprise to him, and he probably thought it would be a surprise to his readers, as well. That's why he offered evidence that the current was still flowing: to show that the battery hadn't just died. The induction effect occurred during the short time it took for the black coil's magnetic field to be established, b/2. Even more counterintuitively, we get an effect, equally strong but in the opposite direction, when the circuit is *broken*, b/4. The effect occurs only when the magnetic field is changing, and it appears to be proportional to the derivative $\partial\mathbf{B}/\partial t$, which is in one direction when the field is being established, and in the opposite direction when it collapses.

The effect is proportional to $\partial\mathbf{B}/\partial t$, but what *is* the effect? A voltmeter is nothing more than a resistor with an attachment for measuring the current through it. A current will not flow through a resistor unless there is some electric field pushing the electrons, so we conclude that the changing *magnetic field* has produced an *electric field* in the surrounding space. Since the white wire is not a perfect conductor, there must be electric fields in it as well. The remarkable thing about the circuit formed by the white wire is that as the electrons travel around and around, they are always being pushed forward by electric fields. This violates the loop rule, which says that when an electron makes a round trip, there is supposed to be just as much "uphill" (moving against the electric field) as "downhill" (moving with it). That's OK. The loop rule is only true for statics. Faraday's experiments show that an electron really can go around and around, and always be going "downhill," as in the famous drawing by M.C. Escher shown in figure c. That's just what happens when you have a curly field.

When a field is curly, we can measure its curliness using a circulation. Unlike the magnetic circulation Γ_B , the electric circulation Γ_E is something we can measure directly using ordinary tools. A circulation is defined by breaking up a loop into tiny segments, ds , and adding up the dot products of these distance vectors with the field. But when we multiply electric field by distance, what we get is an indication of the amount of work per unit charge done on a test charge that has been moved through that distance. The work per unit charge has units of volts, and it can be measured using a voltmeter, as shown in figure e, where Γ_E equals the sum of the voltmeter readings. Since the electric circulation is directly measurable, most people who work with circuits are more familiar with it than they are with the magnetic circulation. They usually refer to Γ_E using the synonym “emf,” which stands for “electromotive force,” and notate it as \mathcal{E} . (This is an unfortunate piece of terminology, because its units are really volts, not newtons.) The term emf can also be used when the path is not a closed loop.

Faraday’s experiment demonstrates a new relationship

$$\Gamma_E \propto -\frac{\partial B}{\partial t},$$

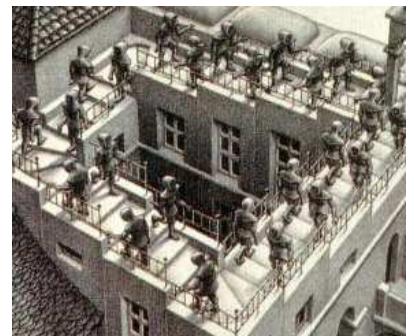
where the negative sign is a way of showing the observed left-handed relationship, d. This is similar to the structure of Ampère’s law:

$$\Gamma_B \propto I_{\text{through}},$$

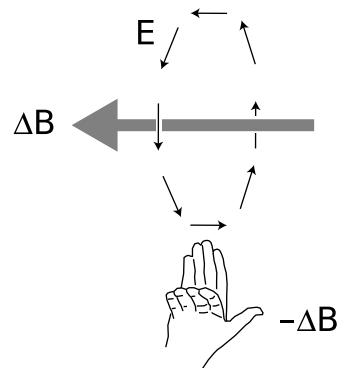
which also relates the curliness of a field to something that is going on nearby (a current, in this case).

It’s important to note that even though the emf, Γ_E , has units of volts, it isn’t a voltage. A voltage is a measure of the electrical energy a charge has when it is at a certain point in space. The curly nature of nonstatic fields means that this whole concept becomes nonsense. In a curly field, suppose one electron stays at home while its friend goes for a drive around the block. When they are reunited, the one that went around the block has picked up some kinetic energy, while the one who stayed at home hasn’t. We simply can’t define an electrical energy $U_e = qV$ so that $U_e + K$ stays the same for each electron. No voltage pattern, V , can do this, because then it would predict the same kinetic energies for the two electrons, which is incorrect. When we’re dealing with nonstatic fields, we need to think of the electrical energy in terms of the energy density of the fields themselves.

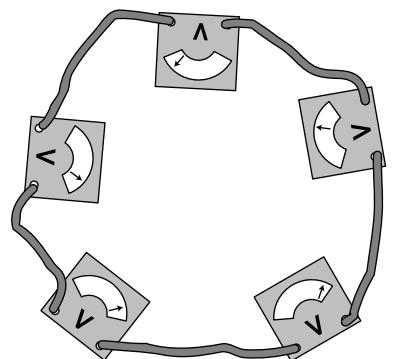
It might sound as though an electron could get a free lunch by circling around and around in a curly electric field, resulting in a violation of conservation of energy. The following examples, in addition to their practical interest, both show that energy is in fact conserved.



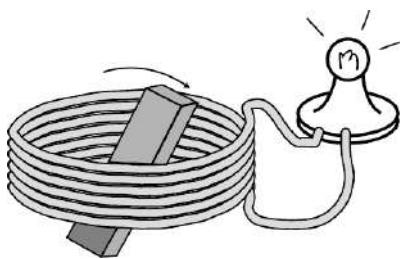
c / Detail from *Ascending and Descending*, M.C. Escher, 1960.



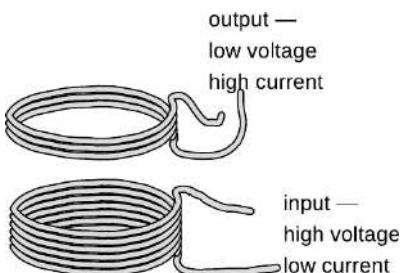
d / The relationship between the change in the magnetic field, and the electric field it produces.



e / The electric circulation is the sum of the voltmeter readings.



f / A generator.



g / A transformer.

The generator

example 18

A basic generator, f, consists of a permanent magnet that rotates within a coil of wire. The magnet is turned by a motor or crank, (not shown). As it spins, the nearby magnetic field changes. This changing magnetic field results in an electric field, which has a curly pattern. This electric field pattern creates a current that whips around the coils of wire, and we can tap this current to light the lightbulb.

If the magnet was on a frictionless bearing, could we light the bulb for free indefinitely, thus violating conservation of energy? No. Mechanical work has to be done to crank the magnet, and that's where the energy comes from. If we break the light-bulb circuit, it suddenly gets easier to crank the magnet! This is because the current in the coil sets up its own magnetic field, and that field exerts a torque on the magnet. If we stopped cranking, this torque would quickly make the magnet stop turning.

self-check G

When you're driving your car, the engine recharges the battery continuously using a device called an alternator, which is really just a generator. Why can't you use the alternator to start the engine if your car's battery is dead?

▷ Answer, p. 1065

The transformer

example 19

In example 18 on page 562, we discussed the advantages of transmitting power over electrical lines using high voltages and low currents. However, we don't want our wall sockets to operate at 10000 volts! For this reason, the electric company uses a device called a transformer, g, to convert everything to lower voltages and higher currents inside your house. The coil on the input side creates a magnetic field. Transformers work with alternating current, so the magnetic field surrounding the input coil is always changing. This induces an electric field, which drives a current around the output coil.

Since the electric field is curly, an electron can keep gaining more and more energy by circling through it again and again. Thus the output voltage can be controlled by changing the number of coils of wire on the output side. Changing the number of coils on the input side also has an effect (homework problem 33).

In any case, conservation of energy guarantees that the amount of power on the output side must equal the amount put in originally, $I_{in}V_{in} = I_{out}V_{out}$, so no matter what factor the voltage is reduced by, the current is increased by the same factor.

Discussion Questions

- A** Suppose the bar magnet in figure f on page 716 has a magnetic field pattern that emerges from its top, circling around and coming back in the bottom. This field is created by electrons orbiting atoms inside

the magnet. Are these atomic currents clockwise or counterclockwise as seen from above? In what direction is the current flowing in the circuit?

We have a circling atomic current inside the circling current in the wires. When we have two circling currents like this, they will make torques on each other that will tend to align them in a certain way. Since currents in the same direction attract one another, which way is the torque made by the wires on the bar magnet? Verify that due to this torque, mechanical work has to be done in order to crank the generator.

11.5.2 Why induction?

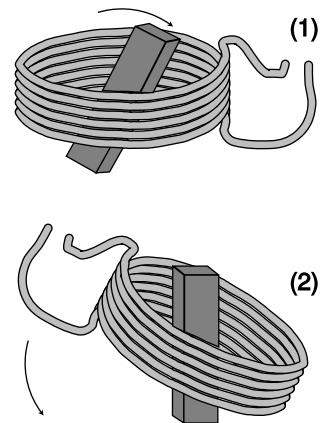
Faraday's results leave us in the dark about several things:

- They don't explain *why* induction effects occur.
- The relationship $\Gamma_E \propto -\partial B / \partial t$ tells us that a changing magnetic field creates an electric field in the surrounding region of space, but the phrase "surrounding region of space" is vague, and needs to be made mathematical.
- Suppose that we can make the "surrounding region of space" idea more well defined. We would then want to know the proportionality constant that has been hidden by the \propto symbol. Although experiments like Faraday's could be used to find a numerical value for this constant, we would like to know why it should have that particular value.

We can get some guidance from the example of a car's alternator (which just means generator), referred to in the self-check on page 716. To keep things conceptually simple, I carefully avoided mentioning that in a real car's alternator, it isn't actually the permanent magnet that spins. The coil is what spins. The choice of design h/1 or h/2 is merely a matter of engineering convenience, not physics. All that matters is the relative motion of the two objects.

This is highly suggestive. As discussed at the beginning of this chapter, magnetism is a relativistic effect. From arguments about relative motion, we concluded that moving electric charges create magnetic fields. Now perhaps we can use reasoning with the same flavor to show that changing magnetic fields produce curly electric fields. Note that figure h/2 doesn't even require induction. The protons and electrons in the coil are moving through a magnetic field, so they experience forces. The protons can't flow, because the coil is a solid substance, but the electrons can, so a current is induced.⁷

Now if we're convinced that figure h/2 produces a current in the coil, then it seems very plausible that the same will happen in



h / It doesn't matter whether it's the coil or the permanent magnet that spins. Either way, we get a functioning generator.

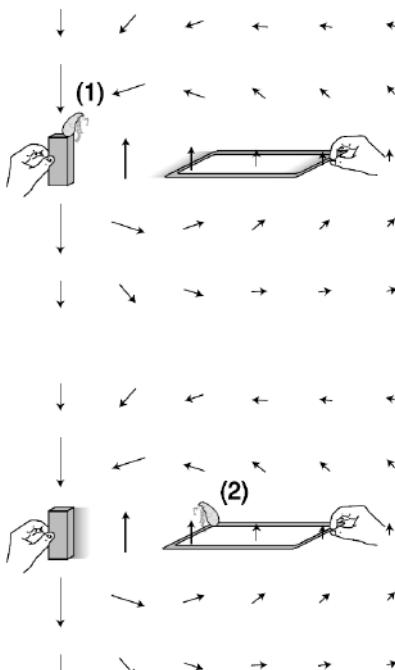
⁷Note that the magnetic field never does work on a charged particle, because its force is perpendicular to the motion; the electric power is actually coming from the mechanical work that had to be done to spin the coil. Spinning the coil is more difficult due to the presence of the magnet.

figure h/1, which implies the existence of induction effects. But this example involves circular motion, so it doesn't quite work as a way of proving that induction exists. When we say that motion is relative, we only mean straight-line motion, not circular motion.

A more ironclad relativistic argument comes from the arrangement shown in figure i. This is also a generator — one that is impractical, but much easier to understand.

Flea 1 doesn't believe in this modern foolishness about induction. She's sitting on the bar magnet, which to her is obviously at rest. As the square wire loop is dragged away from her and the magnet, its protons experience a force out of the page, because the cross product $\mathbf{F} = q\mathbf{v} \times \mathbf{B}$ is out of the page. The electrons, which are negatively charged, feel a force into the page. The conduction electrons are free to move, but the protons aren't. In the front and back sides of the loop, this force is perpendicular to the wire. In the right and left sides, however, the electrons are free to respond to the force. Note that the magnetic field is weaker on the right side. It's as though we had two pumps in a loop of pipe, with the weaker pump trying to push in the opposite direction; the weaker pump loses the argument.⁸ We get a current that circulates around the loop.⁹ There is no induction going on in this frame of reference; the forces that cause the current are just the ordinary magnetic forces experienced by any charged particle moving through a magnetic field.

Flea 2 is sitting on the loop, which she considers to be at rest. In her frame of reference, it's the bar magnet that is moving. Like flea 1, she observes a current circulating around the loop, but unlike flea 1, she cannot use magnetic forces to explain this current. As far as she is concerned, the electrons were initially at rest. Magnetic forces are forces between moving charges and other moving charges, so a magnetic field can never accelerate a charged particle starting from rest. A force that accelerates a charge from rest can only be an *electric* force, so she is forced to conclude that there is an electric field in her region of space. This field drives electrons around and around in circles, so it is apparently violating the loop rule — it is a curly field. What reason can flea 2 offer for the existence of this electric field pattern? Well, she's been noticing that the magnetic



i / A generator that works with linear motion.

⁸If the pump analogy makes you uneasy, consider what would happen if all the electrons moved into the page on both sides of the loop. We'd end up with a net negative charge at the back side, and a net positive charge on the front. This actually would happen in the first nanosecond after the loop was set in motion. This buildup of charge would start to quench both currents due to electrical forces, but the current in the right side of the wire, which is driven by the weaker magnetic field, would be the first to stop. Eventually, an equilibrium will be reached in which the same amount of current is flowing at every point around the loop, and no more charge is being piled up.

⁹The wire is not a perfect conductor, so this current produces heat. The energy required to produce this heat comes from the hands, which are doing mechanical work as they separate the magnet from the loop.

field in her region of space has been changing, possibly because that bar magnet over there has been getting farther away. She observes that a changing magnetic field creates a curly electric field.

We therefore conclude that induction effects *must* exist based on the fact that motion is relative. If we didn't want to admit induction effects, we would have to outlaw flea 2's frame of reference, but the whole idea of relative motion is that all frames of reference are created equal, and there is no way to determine which one is really at rest.

This whole line of reasoning was not available to Faraday and his contemporaries, since they thought the relative nature of motion only applied to matter, not to electric and magnetic fields.¹⁰ But with the advantage of modern hindsight, we can understand in fundamental terms the facts that Faraday had to take simply as mysterious experimental observations. For example, the geometric relationship shown in figure d follows directly from the direction of the current we deduced in the story of the two fleas.

11.5.3 Faraday's law

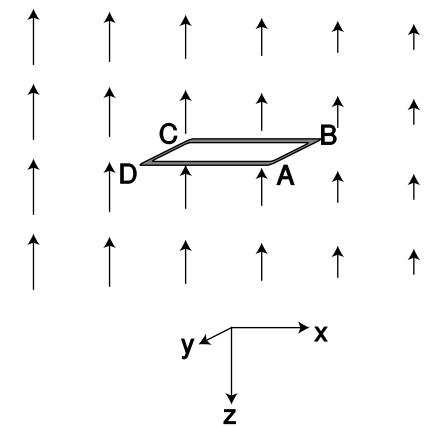
We can also answer the other questions posed on page 717. The divide-and-conquer approach should be familiar by now. We first determine the circulation Γ_E in the case where the wire loop is very tiny, j. Then we can break down any big loop into a grid of small ones; we've already seen that when we make this kind of grid, the circulations add together. Although we'll continue to talk about a physical loop of wire, as in figure i, the tiny loop can really be just like the edges of an Ampèrean surface: a mathematical construct that doesn't necessarily correspond to a real object.

In the close-up view shown in figure j, the field looks simpler. Just as a tiny part of a curve looks straight, a tiny part of this magnetic field looks like the field vectors are just getting shorter by the same amount with each step to the right. Writing dx for the width of the loop, we therefore have

$$B(x + dx) - B(x) = \frac{\partial B}{\partial x} dx$$

for the difference in the strength of the field between the left and right sides. In the frame of reference where the loop is moving, a charge q moving along with the loop at velocity v will experience a magnetic force $\mathbf{F}_B = qvB\hat{\mathbf{y}}$. In the frame moving along with the loop, this is interpreted as an electrical force, $\mathbf{F}_E = qE\hat{\mathbf{y}}$. Observers in the two frames agree on how much force there is, so in the loop's frame, we have an electric field $\mathbf{E} = vB\hat{\mathbf{y}}$. This field is

¹⁰They can't be blamed too much for this. As a consequence of Faraday's work, it soon became apparent that light was an electromagnetic wave, and to reconcile this with the relative nature of motion requires Einstein's version of relativity, with all its subversive ideas how space and time are not absolute.



j / A new version of figure i with a tiny loop. The point of view is above the plane of the loop. In the frame of reference where the magnetic field is constant, the loop is moving to the right.

perpendicular to the front and back sides of the loop, BC and DA, so there is no contribution to the circulation along these sides, but there is a counterclockwise contribution to the circulation on CD, and smaller clockwise one on AB. The result is a circulation that is counterclockwise, and has an absolute value

$$\begin{aligned} |\Gamma_E| &= |E(x) dy - E(x + dx) dy| \\ &= |v[B(x) - B(x + dx)]| dy \\ &= \left| v \frac{\partial B}{\partial x} \right| dx dy \\ &= \left| \frac{dx}{dt} \frac{\partial B}{\partial x} \right| dx dy \\ &= \left| \frac{\partial B}{\partial t} \right| dA. \end{aligned}$$

Using a right-hand rule, the counterclockwise circulation is represented by pointing one's thumb up, but the vector $\partial\mathbf{B}/\partial t$ is down. This is just a rephrasing of the geometric relationship shown in figure d on page 715. We can represent the opposing directions using a minus sign,

$$\Gamma_E = -\frac{\partial B}{\partial t} dA.$$

Although this derivation was carried out with everything aligned in a specific way along the coordinate axes, it turns out that this relationship can be generalized as a vector dot product,

$$\Gamma_E = -\frac{\partial \mathbf{B}}{\partial t} \cdot d\mathbf{A}.$$

Finally, we can take a finite-sized loop and break down the circulation around its edges into a grid of tiny loops. The circulations add, so we have

$$\Gamma_E = -\sum \frac{\partial \mathbf{B}_j}{\partial t} \cdot d\mathbf{A}_j.$$

This is known as Faraday's law. (I don't recommend memorizing all these names.) Mathematically, Faraday's law is very similar to the structure of Ampère's law: the circulation of a field around the edges of a surface is equal to the sum of something that points through the

If the loop itself isn't moving, twisting, or changing shape, then the area vectors don't change over time, and we can move the derivative outside the sum, and rewrite Faraday's law in a slightly more transparent form:

$$\begin{aligned} \Gamma_E &= -\frac{\partial}{\partial t} \sum \mathbf{B}_j \cdot d\mathbf{A}_j \\ &= -\frac{\partial \Phi_B}{\partial t} \end{aligned}$$

A changing magnetic flux makes a curly electric field. You might think based on Gauss' law for magnetic fields that Φ_B would be identically zero. However, Gauss' law only applies to surfaces that are closed, i.e., have no edges.

self-check H

Check that the units in Faraday's law work out. An easy way to approach this is to use the fact that vB has the same units as E , which can be seen by comparing the equations for magnetic and electric forces used above.

▷ Answer, p. 1065

A pathetic generator

example 20

▷ The horizontal component of the earth's magnetic field varies from zero, at a magnetic pole, to about 10^{-4} T near the equator. Since the distance from the equator to a pole is about 10^7 m, we can estimate, very roughly, that the horizontal component of the earth's magnetic field typically varies by about 10^{-11} T/m as you go north or south. Suppose you connect the terminals of a one-ohm lightbulb to each other with a loop of wire having an area of 1 m^2 . Holding the loop so that it lies in the east-west-up-down plane, you run straight north at a speed of 10 m/s, how much current will flow? Next, repeat the same calculation for the surface of a neutron star. The magnetic field on a neutron star is typically 10^9 T, and the radius of an average neutron star is about 10^4 m.

▷ Let's work in the frame of reference of the running person. In this frame of reference, the earth is moving, and therefore the local magnetic field is changing in strength by 10^{-9} T/s. This rate of change is almost exactly the same throughout the interior of the loop, so we can dispense with the summation, and simply write Faraday's law as

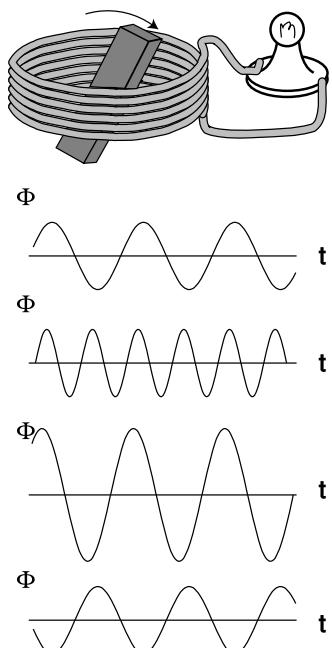
$$\Gamma_E = -\frac{\partial \mathbf{B}}{\partial t} \cdot \mathbf{A}.$$

Since what we estimated was the rate of change of the horizontal component, and the vector \mathbf{A} is horizontal (perpendicular to the loop), we can find this dot product simply by multiplying the two numbers:

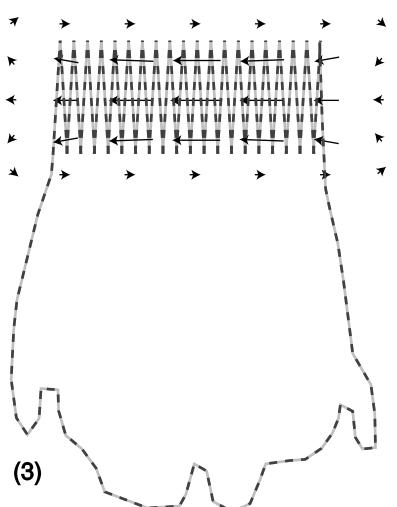
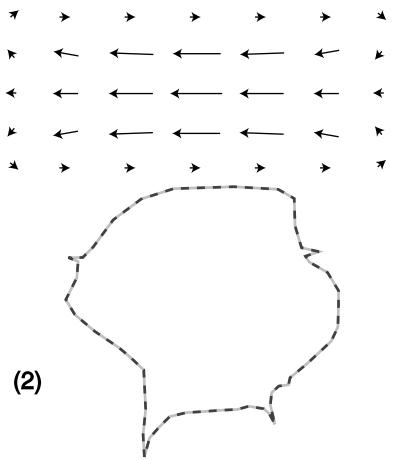
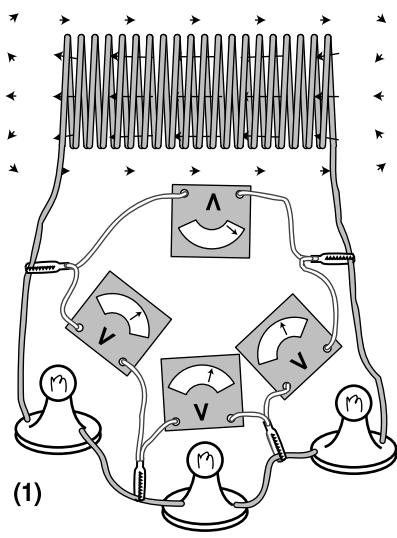
$$\begin{aligned}\Gamma_E &= (10^{-9} \text{ T/s})(1 \text{ m}^2) \\ &= 10^{-9} \text{ T} \cdot \text{m}^2/\text{s} \\ &= 10^{-9} \text{ V}\end{aligned}$$

This is certainly not enough to light the bulb, and would not even be easy to measure using the most sensitive laboratory instruments.

Now what about the neutron star? We'll pretend you're tough enough that its gravity doesn't instantly crush you. The spatial



k / Example 21.



I / Example 22.

variation of the magnetic field is on the order of $(10^9 \text{ T}/10^4 \text{ m}) = 10^5 \text{ T/m}$. If you can run north at the same speed of 10 m/s, then in your frame of reference there is a temporal (time) variation of about 10^6 T/s , and a calculation similar to the previous one results in an emf of 10^6 V ! This isn't just strong enough to light the bulb, it's sufficient to evaporate it, and kill you as well!

It might seem as though having access to a region of rapidly changing magnetic field would therefore give us an infinite supply of free energy. However, the energy that lights the bulb is actually coming from the mechanical work you do by running through the field. A tremendous force would be required to make the wire loop move through the neutron star's field at any significant speed.

Speed and power in a generator

example 21

▷ Figure k shows three graphs of the magnetic flux through a generator's coils as a function of time. In graph 2, the generator is being cranked at twice the frequency. In 3, a permanent magnet with double the strength has been used. In 4, the generator is being cranked in the opposite direction. Compare the power generated in figures 2-4 with the the original case, 1.

▷ If the flux varies as $\Phi = A \sin \omega t$, then the time derivative occurring in Faraday's law is $\partial\Phi/\partial t = A\omega \cos \omega t$. The absolute value of this is the same as the absolute value of the emf, Γ_E . The current through the lightbulb is proportional to this emf, and the power dissipated depends on the square of the current ($P = I^2 R$), so $P \propto A^2 \omega^2$. Figures 2 and 3 both give four times the output power (and require four times the input power). Figure 4 gives the same result as figure 1; we can think of this as a negative amplitude, which gives the same result when squared.

An approximate loop rule

example 22

Figure I/1 shows a simple RL circuit of the type discussed in the last chapter. A current has already been established in the coil, let's say by a battery. The battery was then unclipped from the coil, and we now see the circuit as the magnetic field in and around the inductor is beginning to collapse. I've already cautioned you that the loop rule doesn't apply in nonstatic situations, so we can't assume that the readings on the four voltmeters add up to zero. The interesting thing is that although they don't add up to exactly zero in this circuit, they very nearly do. Why is the loop rule even approximately valid in this situation?

The reason is that the voltmeters are measuring the emf Γ_E around the path shown in figure I/2, and the stray field of the solenoid is extremely weak out there. In the region where the meters are, the arrows representing the magnetic field would be too small to allow me to draw them to scale, so I have simply omitted them. Since the field is so weak in this region, the flux through the loop is nearly zero, and the rate of change of the flux, $\partial\Phi_B/\partial t$, is also

nearly zero. By Faraday's law, then, the emf around this loop is nearly zero.

Now consider figure I/3. The flux through the interior of this path is not zero, because the strong part of the field passes through it, and not just once but many times. To visualize this, imagine that we make a wire frame in this shape, dip it in a tank of soapy water, and pull it out, so that there is a soap-bubble film spanning its interior. Faraday's law refers to the rate of change of the flux through a surface such as this one. (The soap film tends to assume a certain special shape which results in the minimum possible surface area, but Faraday's law would be true for any surface that filled in the loop.) In the coiled part of the wire, the soap makes a three-dimensional screw shape, like the shape you would get if you took the steps of a spiral staircase and smoothed them into a ramp. The loop rule is going to be strongly violated for this path.

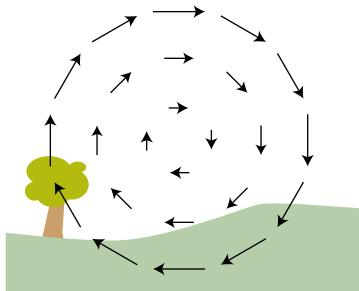
We can interpret this as follows. Since the wire in the solenoid has a very low resistance compared to the resistances of the light bulbs, we can expect that the electric field along the corkscrew part of loop I/3 will be very small. As an electron passes through the coil, the work done on it is therefore essentially zero, and the true emf along the coil is zero. In figure I/1, the meter on top is therefore not telling us the actual emf experienced by an electron that passes through the coil. It is telling us the emf experienced by an electron that passes through the meter itself, which is a different quantity entirely. The other three meters, however, really do tell us the emf through the bulbs, since there are no magnetic fields where they are, and therefore no funny induction effects.

11.6 Maxwell's equations

11.6.1 Induced magnetic fields



a / James Clerk Maxwell (1831–1879)



b / Where is the moving charge responsible for this magnetic field?

We are almost, but not quite, done figuring out the complete set of physical laws, called Maxwell's equations, governing electricity and magnetism. We are only missing one more term. For clarity, I'll state Maxwell's equations with the missing part included, and then discuss the physical motivation and experimental evidence for sticking it in:

Maxwell's equations

For any closed surface, the fluxes through the surface are

$$\Phi_E = 4\pi k q_{in} \quad \text{and}$$
$$\Phi_B = 0.$$

For any surface that is not closed, the circulations around the edges of the surface are given by

$$\Gamma_E = -\frac{\partial \Phi_B}{\partial t} \quad \text{and}$$
$$c^2 \Gamma_B = \frac{\partial \Phi_E}{\partial t} + 4\pi k I_{through}.$$

The Φ_E equation is Gauss' law: charges make diverging electric fields. The corresponding equation for Φ_B tells us that magnetic “charges” (monopoles) don't exist, so magnetic fields never diverge. The third equation says that changing magnetic fields induce curly electric fields, whose curliness we can measure using the emf, Γ_E , around a closed loop. The final equation, for Γ_B , is the only one where anything new has been added. Without the new time derivative term, this equation would simply be Ampère's law. (I've chosen to move the c^2 over to the left because it simplifies the writing, and also because it more clearly demonstrates the analogous roles played by charges and currents in the Φ_E and Γ_B equations.)

This new $\partial \Phi_E / \partial t$ term says that just as a changing magnetic field can induce a curly electric field, a changing electric field can induce a curly magnetic field. Why should this be so? The following examples show that Maxwell's equations would not make sense in general without it.

Figure b shows a mysterious curly magnetic field. Magnetic fields are supposed to be made by moving charges, but there don't seem to be any moving charges in this landscape. Where are they? One reasonable guess would be that they're behind your head, where you can't see them. Suppose there's a positively charged particle about to hit you in the back of the head. This particle is like a current

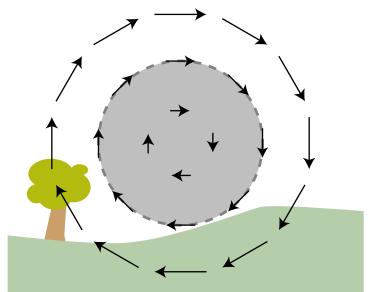
going into the page. We're used to dealing with currents made by many charged particles, but logically we can't have some minimum number that would qualify as a current. This is not a static current, however, because the current at a given point in space is not staying the same over time. If the particle is pointlike, then it takes zero time to pass any particular location, and the current is then infinite at that point in space. A moment later, when the particle is passing by some other location, there will be an infinite current there, and zero current in the previous location. If this single particle qualifies as a current, then it should be surrounded by a curly magnetic field, just like any other current.¹¹

This explanation is simple and reasonable, but how do we know it's correct? Well, it makes another prediction, which is that the positively charged particle should be making an electric field as well. Not only that, but if it's headed for the back of your head, then it's getting closer and closer, so the electric field should be getting stronger over time. But this is exactly what Maxwell's equations require. There is no current I_{through} piercing the Ampèrian surface shown in figure c, so Maxwell's equation for Γ_B becomes $c^2\Gamma_B = \partial\Phi_E/\partial t$. The only reason for an electric field to change is if there are charged particles making it, and those charged particles are moving. When charged particles are moving, they make magnetic fields as well.

Note that the above example is also sufficient to prove the positive sign of the $\partial\Phi_E/\partial t$ term in Maxwell's equations, which is different from the negative sign of Faraday's $-\partial\Phi_B/\partial t$ term.

The addition of the $\partial\Phi_E/\partial t$ term has an even deeper and more important physical meaning. With the inclusion of this term, Maxwell's equations can describe correctly the way in which disturbances in the electric and magnetic fields ripple outwards at the speed of light. Indeed, Maxwell was the first human to understand that light was in fact an electromagnetic wave. Legend has it that he first realized this implication of his equations. He went for a walk with his wife, and told her she was the only other person in the world who really knew what starlight was.

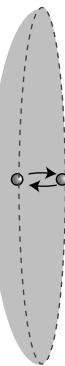
To see how the $\partial\Phi_E/\partial t$ term relates to electromagnetic waves, let's look at an example where we would get nonsense without it. Figure d shows an electron that sits just on one side of an imaginary Ampèrian surface, and then hops through it at some randomly



c / An Ampèrian surface superimposed on the landscape.



d / An electron jumps through a hoop.



e / An alternative Ampèrian surface.

¹¹One way to prove this rigorously is that in a frame of reference where the particle is at rest, it has an electric field that surrounds it on all sides. If the particle has been moving with constant velocity for a long time, then this is just an ordinary Coulomb's-law field, extending off to very large distances, since disturbances in the field ripple outward at the speed of light. In a frame where the particle is moving, this pure electric field is experienced instead as a combination of an electric field and a magnetic field, so the magnetic field must exist throughout the same vast region of space.

chosen moment. Unadorned with the $\partial\Phi_E/\partial t$ term, Maxwell's equation for Γ_B reads as $c^2\Gamma_B = 4\pi kI_{through}$, which is Ampère's law. If the electron is a pointlike particle, then we have an infinite current $I_{through}$ at the moment when it pierces the imaginary surface, and zero current at all other times. An infinite magnetic circulation Γ_B can only be produced by an infinite magnetic field, so without the $\partial\Phi_E/\partial t$ term, Maxwell's equations predict nonsense: the edge of the surface would experience an infinite magnetic field at one instant, and zero magnetic field at all other times. Even if the infinity didn't upset us, it doesn't make sense that anything special would happen at the moment the electron passed through the surface, because the surface is an imaginary mathematical construct. We could just as well have chosen the curved surface shown in figure e, which the electron never crosses at all. We are already clearly getting nonsensical results by omitting the $\partial\Phi_E/\partial t$ term, and this shouldn't surprise us because Ampère's law only applies to statics. More to the point, Ampère's law doesn't have time in it, so it predicts that this effect is instantaneous. According to Ampère's law, we could send Morse code signals by wiggling the electron back and forth, and these signals would be received at distant locations instantly, without any time delay at all. This contradicts the theory of relativity, one of whose predictions is that information cannot be transmitted at speeds greater than the speed of light.

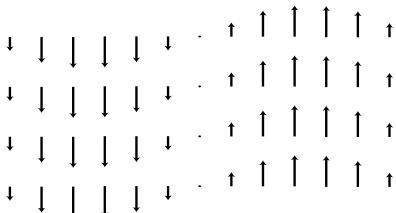
Discussion Questions

A Induced magnetic fields were introduced in the text via the imaginary landscape shown in figure b on page 724, and I argued that the magnetic field could have been produced by a positive charge coming from behind your head. This is a specific assumption about the *number* of charges (one), the *direction* of motion, and the *sign* of the charge. What are some other scenarios that could explain this field?

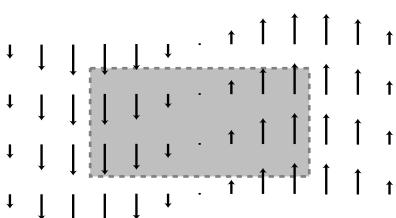
11.6.2 Light waves

We could indeed send signals using this scheme, and the signals would be a form of light. A radio transmitting antenna, for instance, is simply a device for whipping electrons back and forth at megahertz frequencies. Radio waves are just like visible light, but with a lower frequency. With the addition of the $\partial\Phi_E/\partial t$ term, Maxwell's equations are capable of describing electromagnetic waves. It would be possible to use Maxwell's equations to calculate the pattern of the electric and magnetic fields rippling outward from a single electron that fidgets at irregular intervals, but let's pick a simpler example to analyze.

The simplest wave pattern is a sine wave like the one shown in figure f. Let's assume a magnetic field of this form, and see what Maxwell's equations tell us about it. If the wave is traveling through empty space, then there are no charges or currents present,



f / A magnetic field in the form of a sine wave.



g / The wave pattern is curly. For example, the circulation around this rectangle is nonzero and counterclockwise.

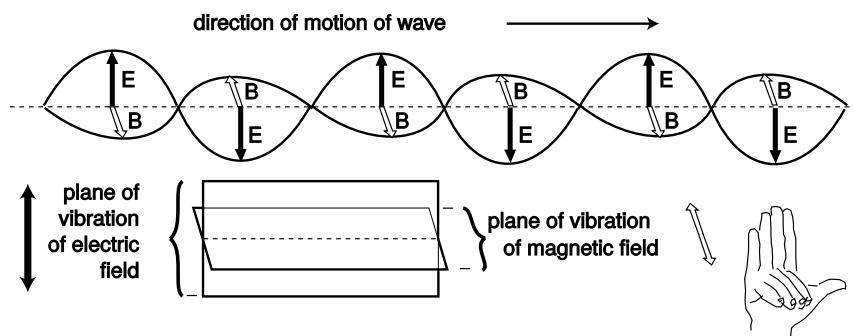
and Maxwell's equations become

$$\begin{aligned}\Phi_E &= 0 \\ \Phi_B &= 0 \\ \Gamma_E &= -\frac{\partial \Phi_B}{\partial t} \\ c^2 \Gamma_B &= \frac{\partial \Phi_E}{\partial t}.\end{aligned}$$

The equation $\Phi = 0$ has already been verified for this type of wave pattern in example 39 on page 654. Even if you haven't learned the techniques from that section, it should be visually plausible that this field pattern doesn't diverge or converge on any particular point.

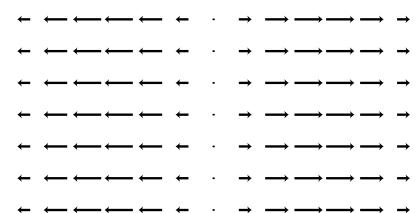
Geometry of the electric and magnetic fields

The equation $c^2 \Gamma_B = \partial \Phi_E / \partial t$ tells us that there can be no such thing as a purely magnetic wave. The wave pattern clearly does have a nonvanishing circulation around the edge of the surface suggested in figure g, so there must be an electric flux through the surface. This magnetic field pattern must be intertwined with an electric field pattern that fills the same space. There is also no way that the two sides of the equation could stay synchronized with each other unless the electric field pattern is also a sine wave, and one that has the same wavelength, frequency, and velocity. Since the electric field is making a flux through the indicated surface, it's plausible that the electric field vectors lie in a plane perpendicular to that of the magnetic field vectors. The resulting geometry is shown in figure h. Further justification for this geometry is given later in this subsection.



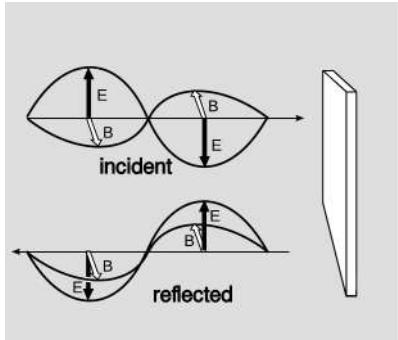
h / The geometry of an electromagnetic wave.

One feature of figure h that is easily justified is that the electric and magnetic fields are perpendicular not only to each other, but also to the direction of propagation of the wave. In other words, the vibration is sideways, like people in a stadium "doing the wave," not lengthwise, like the accordion pattern in figure i. (In standard



i / An impossible wave pattern.

wave terminology, we say that the wave is transverse, not longitudinal.) The wave pattern in figure i is impossible, because it diverges from the middle. For virtually any choice of Gaussian surface, the magnetic and electric fluxes would be nonzero, contradicting the equations $\Phi_B = 0$ and $\Phi_E = 0$.¹²



j / Example 23. The incident and reflected waves are drawn offset from each other for clarity, but are actually on top of each other so that their fields superpose.

Reflection

example 23

The wave in figure j hits a silvered mirror. The metal is a good conductor, so it has constant voltage throughout, and the electric field equals zero inside it: the wave doesn't penetrate and is 100% reflected. If the electric field is to be zero at the surface as well, the reflected wave must have its electric field inverted (p. 376), so that the incident and reflected fields cancel there.

But the magnetic field of the reflected wave is *not* inverted. This is because the reflected wave has to have the correct right-handed relationship between the fields and the direction of propagation.

Polarization

Two electromagnetic waves traveling in the same direction through space can differ by having their electric and magnetic fields in different directions, a property of the wave called its polarization.

The speed of light

What is the velocity of the waves described by Maxwell's equations? Maxwell convinced himself that light was an electromagnetic wave partly because his equations predicted waves moving at the velocity of light, c . The only velocity that appears in the equations is c , so this is fairly plausible, although a real calculation is required in order to prove that the velocity of the waves isn't something like $2c$ or c/π — or zero, which is also c multiplied by a constant! The following discussion, leading up to a proof that electromagnetic waves travel at c , is meant to be understandable even if you're reading this book out of order, and haven't yet learned much about waves. As always with proofs in this book, the reason to read it isn't to convince yourself that it's true, but rather to build your intuition. The style will be visual. In all the following figures, the wave patterns are moving across the page (let's say to the right), and it usually doesn't matter whether you imagine them as representing the wave's magnetic field or its electric field, because Maxwell's equations in a vacuum have the same form for both fields. Whichever field we imagine the figures as representing, the other field is coming in and out of the page.

The velocity of the waves is not zero. If the wave pattern was

¹²Even if the fields can't be parallel to the direction of propagation, one might wonder whether they could form some angle other than 90 degrees with it. No. One proof is given on page 732. An alternative argument, which is simpler but more esoteric, is that if there was such a pattern, then there would be some other frame of reference in which it would look like figure i.

standing still in space, then the right sides of the Γ equations would be zero, because there would be no change in the field over time at a particular point. But the left sides are not zero, so this is impossible.¹³

The velocity of the waves is a fixed number for a given wave pattern. Consider a typical sinusoidal wave of visible light, with a distance of half a micrometer from one peak to the next peak. Suppose this wave pattern provides a valid solution to Maxwell's equations when it is moving with a certain velocity. We then know, for instance, that there *cannot* be a valid solution to Maxwell's equations in which the same wave pattern moves with double that velocity. The time derivatives on the right sides of Maxwell's equations for Γ_E and Γ_B would be twice as big, since an observer at a certain point in space would see the wave pattern sweeping past at twice the rate. But the left sides would be the same, so the equations wouldn't equate.

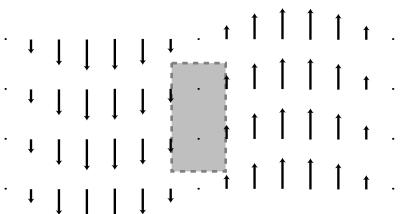
The velocity is the same for all wave patterns. In other words, it isn't $0.878c$ for one wave pattern, and $1.067c$ for some other pattern. This is surprising, since, for example, water waves with different shapes do travel at different speeds. Similarly, even though we speak of "the speed of sound," sound waves do travel at slightly different speeds depending on their pitch and loudness, although the differences are small unless you're talking about cannon blasts or extremely high frequency ultrasound. To see how Maxwell's equations give a consistent velocity, consider figure k. Along the right and left edges of the same Ampèrean surface, the more compressed wave pattern of blue light has twice as strong a field, so the circulations on the left sides of Maxwell's equations are twice as large.¹⁴ To satisfy Maxwell's equations, the time derivatives of the fields must also be twice as large for the blue light. But this is true only if the blue light's wave pattern is moving to the right at the *same* speed as the red light's: if the blue light pattern is sweeping over an observer with a given velocity, then the time between peaks is half as much, like the clicking of the wheels on a train whose cars are half the length.¹⁵

We can also check that bright and dim light, as shown in figure l, have the same velocity. If you haven't yet learned much about

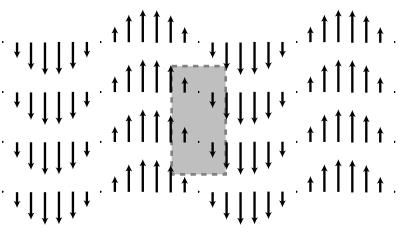
¹³A young Einstein worried about what would happen if you rode a motorcycle alongside a light wave, traveling at the speed of light. Would the light wave have a zero velocity in this frame of reference? The only solution lies in the theory of relativity, one of whose consequences is that a material object like a student or a motorcycle cannot move at the speed of light.

¹⁴Actually, this is only exactly true if the rectangular strip is made infinitesimally thin.

¹⁵You may know already that different colors of light have different speeds when they pass through a material substance, such as the glass or water. This is not in contradiction with what I'm saying here, since this whole analysis is for light in a vacuum.

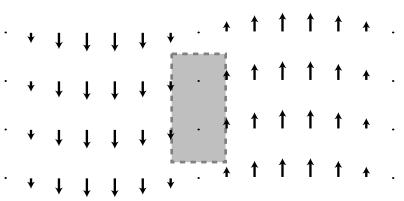


red light

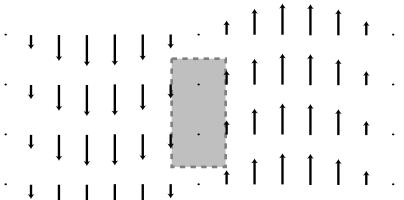


blue light

k / Red and blue light travel at the same speed.

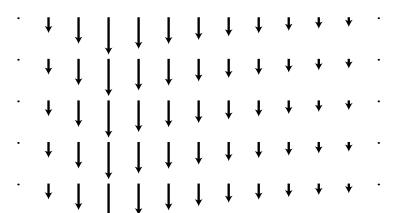


dim light



bright light

l / Bright and dim light travel at the same speed.



m / A nonsinusoidal wave.

waves, then this might be surprising. A material object with more energy goes faster, but that's not the case for waves. The circulation around the edge of the Ampèrean surface shown in the figure is twice as strong for the light whose fields are doubled in strength, so the left sides of Maxwell's Γ equations are doubled. The right sides are also doubled, because the derivative of twice a function is twice the derivative of the original function. Thus if dim light moving with a particular velocity is a solution, then so is bright light, provided that it has the same velocity.

We can now see that all sinusoidal waves have the same velocity. What about nonsinusoidal waves like the one in figure m? There is a mathematical theorem, due to Fourier, that says any function can be made by adding together sinusoidal functions. For instance, $3\sin x - 7\cos 3x$ can be made by adding together the functions $3\sin x$ and $-7\cos 3x$, but Fourier proved that this can be done even for functions, like figure m, that aren't obviously built out of sines and cosines in the first place. Therefore our proof that sinusoidal waves all have the same velocity is sufficient to demonstrate that other waves also have this same velocity.

We're now ready to prove that this universal speed for all electromagnetic waves is indeed c . Since we've already convinced ourselves that all such waves travel at the same speed, it's sufficient to find the velocity of one wave in particular. Let's pick the wave whose fields have magnitudes

$$E = \tilde{E} \sin(x + vt) \quad \text{and} \\ B = \tilde{B} \sin(x + vt),$$

which is about as simple as we can get. The peak electric field of this wave has a strength \tilde{E} , and the peak magnetic field is \tilde{B} . The sine functions go through one complete cycle as x increases by $2\pi = 6.28\dots$, so the distance from one peak of this wave to the next — its wavelength — is 6.28... meters. This means that it is not a wave of visible light but rather a radio wave (its wavelength is on the same order of magnitude as the size of a radio antenna). That's OK. What was glorious about Maxwell's work was that it unified the whole electromagnetic spectrum. Light is simple. Radio waves aren't fundamentally any different than light waves, x-rays, or gamma rays.¹⁶

The justification for putting $x + vt$ inside the sine functions is as follows. As the wave travels through space, the whole pattern just shifts over. The fields are zero at $x = 0, t = 0$, since the sine of zero is zero. This zero-point of the wave pattern shifts over as time goes by; at any time t its location is given by $x + vt = 0$. After one

¹⁶What makes them appear to be unrelated phenomena is that we experience them through their interaction with atoms, and atoms are complicated, so they respond to various kinds of electromagnetic waves in complicated ways.

second, the zero-point is located at $x = -(1 \text{ s})v$. The distance it travels in one second is therefore numerically equal to v , and this is exactly the concept of velocity: how far something goes per unit time.

The wave has to satisfy Maxwell's equations for Γ_E and Γ_B regardless of what Ampèrean surfaces we pick, and by applying them to any surface, we could determine the speed of the wave. The surface shown in figure n turns out to result in an easy calculation: a narrow strip of width 2ℓ and height h , coinciding with the position of the zero-point of the field at $t = 0$.

Now let's apply the equation $c^2\Gamma_B = \partial\Phi_E/\partial t$ at $t = 0$. Since the strip is narrow, we can approximate the magnetic field using $\sin x \approx x$, which is valid for small x . The magnetic field on the right edge of the strip, at $x = \ell$, is then $\tilde{B}\ell$, so the right edge of the strip contributes $\tilde{B}\ell h$ to the circulation. The left edge contributes the same amount, so the left side of Maxwell's equation is

$$c^2\Gamma_B = c^2 \cdot 2\tilde{B}\ell h.$$

The other side of the equation is

$$\begin{aligned} \frac{\partial\Phi_E}{\partial t} &= \frac{\partial}{\partial t}(EA) \\ &= 2\ell h \frac{\partial E}{\partial t}, \end{aligned}$$

where we can dispense with the usual sum because the strip is narrow and there is no variation in the field as we go up and down the strip. The derivative equals $v\tilde{E} \cos(x + vt)$, and evaluating the cosine at $x = 0$, $t = 0$ gives

$$\frac{\partial\Phi_E}{\partial t} = 2v\tilde{E}\ell h$$

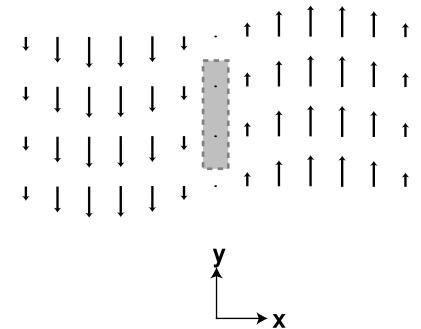
Maxwell's equation for Γ_B therefore results in

$$\begin{aligned} 2c^2\tilde{B}\ell h &= 2\tilde{E}\ell hv \\ c^2\tilde{B} &= v\tilde{E}. \end{aligned}$$

An application of $\Gamma_E = -\partial\Phi_B/\partial t$ gives a similar result, except that there is no factor of c^2

$$\tilde{E} = v\tilde{B}.$$

(The minus sign simply represents the right-handed relationship of the fields relative to their direction of propagation.)



n / The magnetic field of the wave. The electric field, not shown, is perpendicular to the page.

Multiplying these last two equations by each other, we get

$$\begin{aligned}c^2 \tilde{B} \tilde{E} &= v^2 \tilde{E} \tilde{B} \\c^2 &= v^2 \\v &= \pm c.\end{aligned}$$

This is the desired result. (The plus or minus sign shows that the wave can travel in either direction.)

As a byproduct of this calculation, we can find the relationship between the strengths of the electric and magnetic fields in an electromagnetic wave. If, instead of multiplying the equations $c^2 \tilde{B} = v \tilde{E}$ and $\tilde{E} = v \tilde{B}$, we divide them, we can easily show that $\tilde{E} = c \tilde{B}$.

o / The electromagnetic spectrum.

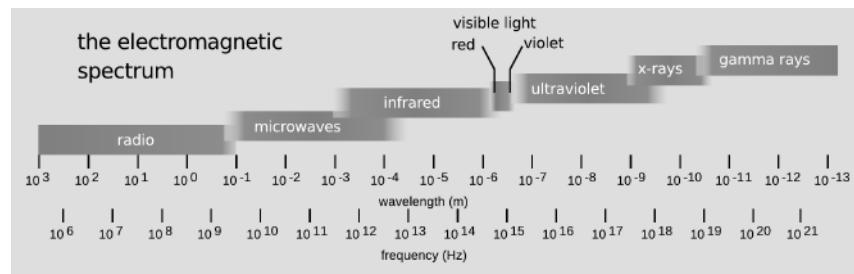


Figure o shows the complete spectrum of light waves. The wavelength λ (number of meters per cycle) and frequency f (number of cycles per second) are related by the equation $c = f\lambda$. Maxwell's equations predict that all light waves have the same structure, regardless of wavelength and frequency, so even though radio and x-rays, for example, hadn't been discovered, Maxwell predicted that such waves would have to exist. Maxwell's 1865 prediction passed an important test in 1888, when Heinrich Hertz published the results of experiments in which he showed that radio waves could be manipulated in the same ways as visible light waves. Hertz showed, for example, that radio waves could be reflected from a flat surface, and that the directions of the reflected and incoming waves were related in the same way as with light waves, forming equal angles with the surface. Likewise, light waves can be focused with a curved, dish-shaped mirror, and Hertz demonstrated the same thing with radio waves using a metal dish.

Momentum of light waves

A light wave consists of electric and magnetic fields, and fields contain energy. Thus a light wave carries energy with it when it travels from one place to another. If a material object has kinetic energy and moves from one place to another, it must also have momentum, so it is logical to ask whether light waves have momentum as well. It can be proved based on relativity¹⁷ that it does, and that the

¹⁷See problem 11 on p. 460, or example 24 on p. 438.

momentum and energy are related by the equation $U = p/c$, where p is the magnitude of the momentum vector, and $U = U_e + U_m$ is the sum of the energy of the electric and magnetic fields. We can now demonstrate this without explicitly referring to relativity, and connect it to the specific structure of a light wave.

The energy density of a light wave is related to the magnitudes of the fields in a specific way — it depends on the squares of their magnitudes, E^2 and B^2 , which are the same as the dot products $\mathbf{E} \cdot \mathbf{E}$ and $\mathbf{B} \cdot \mathbf{B}$. We argued on page 606 that since energy is a scalar, the only possible expressions for the energy densities of the fields are dot products like these, multiplied by some constants. This is because the dot product is the only mathematically sensible way of multiplying two vectors to get a scalar result. (Any other way violates the symmetry of space itself.)

How does this relate to momentum? Well, we know that if we double the strengths of the fields in a light beam, it will have four times the energy, because the energy depends on the square of the fields. But we then know that this quadruple-energy light beam must have quadruple the momentum as well. If there wasn't this kind of consistency between the momentum and the energy, then we could violate conservation of momentum by combining light beams or splitting them up. We therefore know that the momentum density of a light beam must depend on a field multiplied by a field. Momentum, however, is a vector, and there is only one physically meaningful way of multiplying two vectors to get a vector result, which is the cross product (see page 1027). The momentum density can therefore only depend on the cross products $\mathbf{E} \times \mathbf{E}$, $\mathbf{B} \times \mathbf{B}$, and $\mathbf{E} \times \mathbf{B}$. But the first two of these are zero, since the cross product vanishes when there is a zero angle between the vectors. Thus the momentum per unit volume must equal $\mathbf{E} \times \mathbf{B}$ multiplied by some constant,

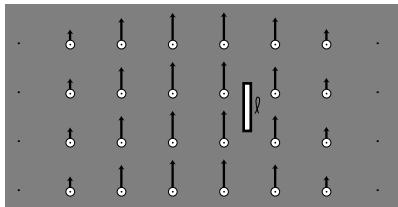
$$d\mathbf{p} = (\text{constant}) \mathbf{E} \times \mathbf{B} dv$$

This predicts something specific about the direction of propagation of a light wave: it must be along the line perpendicular to the electric and magnetic fields. We've already seen that this is correct, and also that the electric and magnetic fields are perpendicular to each other. Therefore this cross product has a magnitude

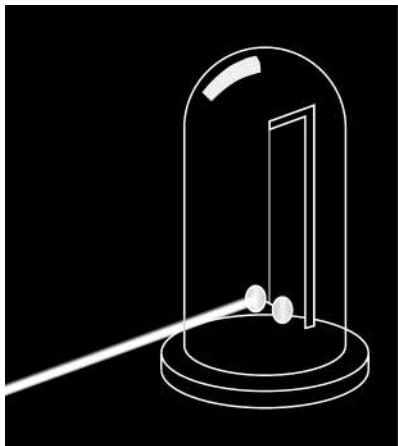
$$\begin{aligned} |\mathbf{E} \times \mathbf{B}| &= |\mathbf{E}| |\mathbf{B}| \sin 90^\circ \\ &= |\mathbf{E}| |\mathbf{B}| \\ &= \frac{|\mathbf{E}|^2}{c} = c |\mathbf{B}|^2, \end{aligned}$$

where in the last step the relation $|\mathbf{E}| = c |\mathbf{B}|$ has been used.

We now only need to find one physical example in order to fix the constant of proportionality. Indeed, if we didn't know relativity, it would be possible to believe that the constant of proportionality was



p / A classical calculation of the momentum of a light wave. An antenna of length ℓ is bathed in an electromagnetic wave. The black arrows represent the electric field, the white circles the magnetic field coming out of the page. The wave is traveling to the right.



q / A simplified drawing of the 1903 experiment by Nichols and Hull that verified the predicted momentum of light waves. Two circular mirrors were hung from a fine quartz fiber, inside an evacuated bell jar. A 150 mW beam of light was shone on one of the mirrors for 6 s, producing a tiny rotation, which was measurable by an optical lever (not shown). The force was within 0.6% of the theoretically predicted value (problem 11 on p. 460) of 0.001 μN . For comparison, a short clipping of a single human hair weighs $\sim 1 \mu\text{N}$.

zero! The simplest example of which I know is as follows. Suppose a piece of wire of length ℓ is bathed in electromagnetic waves coming in sideways, and let's say for convenience that this is a radio wave, with a wavelength that is large compared to ℓ , so that the fields don't change significantly across the length of the wire. Let's say the electric field of the wave happens to be aligned with the wire. Then there is an emf between the ends of the wire which equals $E\ell$, and since the wire is small compared to the wavelength, we can pretend that the field is uniform, not curly, in which case voltage is a well-defined concept, and this is equivalent to a voltage difference $\Delta V = E\ell$ between the ends of the wire. The wire obeys Ohm's law, and a current flows in response to the wave.¹⁸ Equating the expressions dU/dt and $I\Delta V$ for the power dissipated by ohmic heating, we have

$$dU = IE\ell dt$$

for the energy the wave transfers to the wire in a time interval dt .

Note that although some electrons have been set in motion in the wire, we haven't yet seen any momentum transfer, since the protons are experiencing the same amount of electric force in the opposite direction. However, the electromagnetic wave also has a magnetic field, and a magnetic field transfers momentum to (exerts a force on) a current. This is only a force on the electrons, because they're what make the current. The magnitude of this force equals ℓIB (homework problem 6), and using the definition of force, $d\mathbf{p}/dt$, we find for the magnitude of the momentum transferred:

$$dp = \ell IB dt$$

We now know both the amount of energy and the amount of momentum that the wave has lost by interacting with the wire. Dividing these two equations, we find

$$\begin{aligned} \frac{dp}{dU} &= \frac{B}{E} \\ &= \frac{1}{c}, \end{aligned}$$

which is what we expected based on relativity. This can now be restated in the form $d\mathbf{p} = (\text{constant})\mathbf{E} \times \mathbf{B} dv$ (homework problem 40).

Note that although the equations $p = U/c$ and $d\mathbf{p} = (\text{constant})\mathbf{E} \times \mathbf{B} dv$ are consistent with each other for a sine wave, they are not consistent with each other in general. The relativistic argument leading up to $p = U/c$ assumed that we were only talking about

¹⁸This current will soon come to a grinding halt, because we don't have a complete circuit, but let's say we're talking about the first picosecond during which the radio wave encounters the wire.

a single thing traveling in a single direction, whereas no such assumption was made in arguing for the $\mathbf{E} \times \mathbf{B}$ form. For instance, if two light beams of equal strength are traveling through one another, going in opposite directions, their total momentum is zero, which is consistent with the $\mathbf{E} \times \mathbf{B}$ form, but not with U/c .

Some examples were given in chapter 3 of situations where it actually matters that light has momentum. Figure q shows the first confirmation of this fact in the laboratory.

Angular momentum of light waves

For completeness, we note that since light carries momentum, it must also be possible for it to have angular momentum. If you've studied chemistry, here's an example of why this can be important. You know that electrons in atoms can exist in states labeled s, p, d, f, and so on. What you might not have realized is that these are angular momentum labels. The s state, for example, has zero angular momentum. If light didn't have angular momentum, then, for example, it wouldn't be possible for a hydrogen atom in a p state to change to the lower-energy s state by emitting light. Conservation of angular momentum requires that the light wave carry away all the angular momentum originally possessed by the electron in the p state, since in the s state it has none.

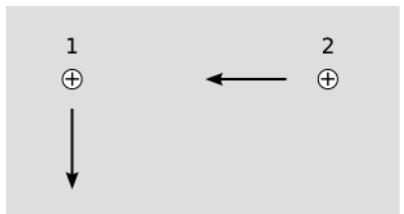
Discussion Questions

A Positive charges 1 and 2 are moving as shown. What electric and magnetic forces do they exert on each other? What does this imply for conservation of momentum?

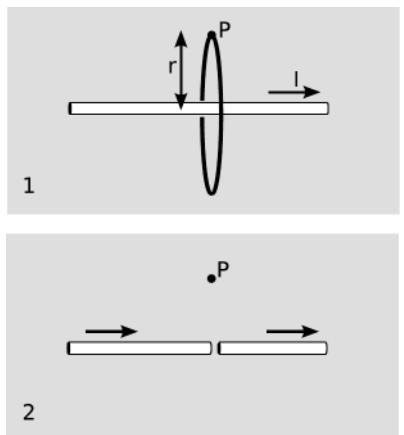
B 1. The figure shows a line of charges moving to the right, creating a current I . An Ampèrean surface in the form of a disk has been superimposed. Use Maxwell's equations to find the field B at point P.
2. A tiny gap is chopped out of the line of charge. What happens when this gap is directly underneath the point P?

C The diagram shows an electric field pattern frozen at one moment in time. Let's imagine that it's the electric part of an electromagnetic wave. Consider four possible directions in which it could be propagating: left, right, up, and down. Determine whether each of these is consistent with Maxwell's equations. If so, infer the direction of the magnetic field.

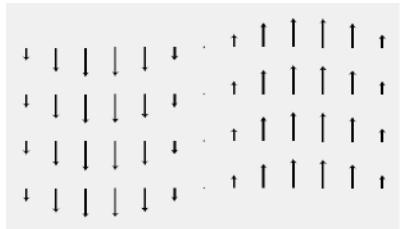
D What happens if we use Maxwell's equations to analyze the behavior of the wave in a frame of reference moving along with the wave?



Discussion question A.



Discussion question B.



Discussion questions C and D.

11.7 Electromagnetic properties of materials

Different types of matter have a variety of useful electrical and magnetic properties. Some are conductors, and some are insulators. Some, like iron and nickel, can be magnetized, while others have useful electrical properties, e.g., dielectrics, discussed qualitatively in the discussion question on page 616, which allow us to make capacitors with much higher values of capacitance than would otherwise be possible. We need to organize our knowledge about the properties that materials can possess, and see whether this knowledge allows us to calculate anything useful with Maxwell's equations.

11.7.1 Conductors

A perfect conductor, such as a superconductor, has no DC electrical resistance. It is not possible to have a static electric field inside it, because then charges would move in response to that field, and the motion of the charges would tend to reduce the field, contrary to the assumption that the field was static. Things are a little different at the surface of a perfect conductor than on the interior. We expect that any net charges that exist on the conductor will spread out under the influence of their mutual repulsion, and settle on the surface. As we saw in chapter 10, Gauss's law requires that the fields on the two sides of a sheet of charge have $|\mathbf{E}_{\perp,1} - \mathbf{E}_{\perp,2}|$ proportional to the surface charge density, and since the field inside the conductor is zero, we infer that there can be a field on or immediately outside the conductor, with a nonvanishing component perpendicular to the surface. The component of the field parallel to the surface must vanish, however, since otherwise it would cause the charges to move along the surface.

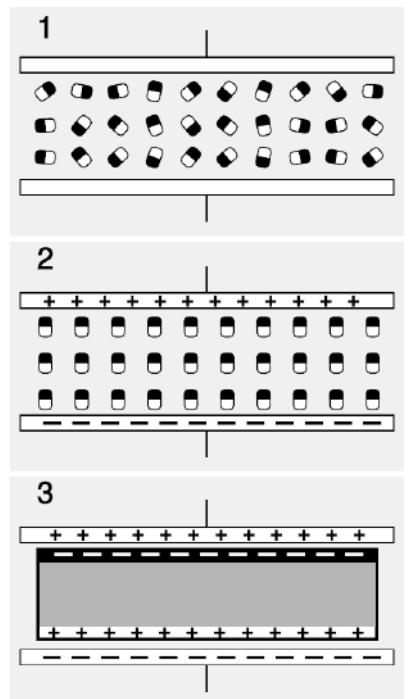
On a hot summer day, the reason the sun feels warm on your skin is that the oscillating fields of the light waves excite currents in your skin, and these currents dissipate energy by ohmic heating. In a perfect conductor, however, this could never happen, because there is no such thing as ohmic heating. Since electric fields can't penetrate a perfect conductor, we also know that an electromagnetic wave can never pass into one. By conservation of energy, we know that the wave can't just vanish, and if the energy can't be dissipated as heat, then the only remaining possibility is that all of the wave's energy is reflected. This is why metals, which are good electrical conductors, are also highly reflective. They are not *perfect* electrical conductors, however, so they are not perfectly reflective. The wave enters the conductor, but immediately excites oscillating currents, and these oscillating currents dissipate the energy both by ohmic heating and by reradiating the reflected wave. Since the parts of Maxwell's equations describing radiation have time derivatives in them, the efficiency of this reradiation process depends strongly on frequency. When the frequency is high and the material is a good conductor, reflection predominates, and is so efficient that the wave

only penetrates to a very small depth, called the skin depth. In the limit of poor conduction and low frequencies, absorption predominates, and the skin depth becomes much greater. In a high-frequency AC circuit, the skin depth in a copper wire is very small, and therefore the signals in such a circuit are propagated entirely at the surfaces of the wires. In the limit of low frequencies, i.e., DC, the skin depth approaches infinity, so currents are carried uniformly over the wires' cross-sections.

We can quantify how well a particular material conducts electricity. We know that the resistance of a wire is proportional to its length, and inversely proportional to its cross-sectional area. The constant of proportionality is $1/\sigma$, where σ (not the same σ as the surface charge density) is called the electrical conductivity. Exposed to an electric field \mathbf{E} , a conductor responds with a current per unit cross-sectional area $\mathbf{J} = \sigma\mathbf{E}$. The skin depth is proportional to $1/\sqrt{f\sigma}$, where f is the frequency of the wave.

11.7.2 Dielectrics

A material with a very low conductivity is an insulator. Such materials are usually composed of atoms or molecules whose electrons are strongly bound to them; since the atoms or molecules have zero total charge, their motion cannot create an electric current. But even though they have zero charge, they may not have zero dipole moment. Imagine such a substance filling in the space between the plates of a capacitor, as in figure a. For simplicity, we assume that the molecules are oriented randomly at first, a/1, and then become completely aligned when a field is applied, a/2. The effect has been to take all of the negatively charged black ends of the molecules and shift them upward, and the opposite for the positively charged white ends. Where the black and white charges overlap, there is still zero net charge, but we have a strip of negative charge at the top, and a strip of positive charge at the bottom, a/3. The effect has been to cancel out part of the charge that was deposited on the plates of the capacitor. Now this is very subtle, because Maxwell's equations treat these charges on an equal basis, but in terms of practical measurements, they are completely different. The charge on the plates can be measured by inserting an ammeter in the circuit, and integrating the current over time. But the charges in the layers at the top and bottom of the dielectric never flowed through any wires, and cannot be detected by an ammeter. In other words, the total charge, q , appearing in Maxwell's equations is actually $q = q_{\text{free}} - q_{\text{bound}}$, where q_{free} is the charge that moves freely through wires, and can be detected in an ammeter, while q_{bound} is the charge bound onto the individual molecules, which can't. We will, however, detect the presence of the bound charges via their electric fields. Since their electric fields partially cancel the fields of the free charges, a voltmeter will register a smaller than expected



a / A capacitor with a dielectric between the plates.

voltage difference between the plates. If we measure q_{free}/V , we have a result that is larger than the capacitance we would have expected.

Although the relationship $\mathbf{E} \leftrightarrow q$ between electric fields and their sources is unalterably locked in by Gauss's law, that's not what we see in practical measurements. In this example, we can measure the voltage difference between the plates of the capacitor and divide by the distance between them to find \mathbf{E} , and then integrate an ammeter reading to find q_{free} , and we will find that Gauss's law appears not to hold. We have $\mathbf{E} \leftrightarrow q_{\text{free}}/(\text{constant})$, where the constant fudge factor is greater than one. This constant is a property of the dielectric material, and tells us how many dipoles there are, how strong they are, and how easily they can be reoriented. The conventional notation is to incorporate this fudge factor into Gauss's law by defining an altered version of the electric field,

$$\mathbf{D} = \epsilon \mathbf{E},$$

and to rewrite Gauss's law as

$$\Phi_D = q_{\text{in, free}}.$$

The constant ϵ is a property of the material, known as its permittivity. In a vacuum, ϵ takes on a value known as ϵ_0 , defined as $1/(4\pi k)$. In a dielectric, ϵ is greater than ϵ_0 . When a dielectric is present between the plates of a capacitor, its capacitance is proportional to ϵ (problem 38). The following table gives some sample values of the permittivities of a few substances.

substance	ϵ/ϵ_0 at zero frequency
vacuum	1
air	1.00054
water	80
barium titanate	1250



b / A stud finder is used to locate the wooden beams, or studs, that form the frame behind the wallboard. It is a capacitor whose capacitance changes when it is brought close to a substance with a particular permittivity. Although the wall is external to the capacitor, a change in capacitance is still observed, because the capacitor has "fringing fields" that extend outside the region between its plates.

A capacitor with a very high capacitance is potentially a superior replacement for a battery, but until the 1990's this was impractical because capacitors with high enough values couldn't be made, even with dielectrics having the largest known permittivities. Such supercapacitors, some with values in the kilofarad range, are now available. Most of them do not use dielectric at all; the very high capacitance values are instead obtained by using electrodes that are not parallel metal plates at all, but exotic materials such as aerogels, which allows the spacing between the "electrodes" to be very small.

Although figure a/2 shows the dipoles in the dielectric being completely aligned, this is not a situation commonly encountered in practice. In such a situation, the material would be as polarized as it could possibly be, and if the field was increased further, it would not respond. In reality, a capacitor, for example, would normally be operated with fields that produced quite a small amount of alignment, and it would be under these conditions that the linear

relationship $\mathbf{D} = \epsilon \mathbf{E}$ would actually be a good approximation. Before a material's maximum polarization is reached, it may actually spark or burn up.

self-check 1

Suppose a parallel-plate capacitor is built so that a slab of dielectric material can be slid in or out. (This is similar to the way the stud finder in figure b works.) We insert the dielectric, hook the capacitor up to a battery to charge it, and then use an ammeter and a voltmeter to observe what happens when the dielectric is withdrawn. Predict the changes observed on the meters, and correlate them with the expected change in capacitance. Discuss the energy transformations involved, and determine whether positive or negative work is done in removing the dielectric.

▷ Answer, p. 1065

11.7.3 Magnetic materials

Magnetic permeability

Atoms and molecules may have magnetic dipole moments as well as electric dipole moments. Just as an electric dipole contains bound charges, a magnetic dipole has bound currents, which come from the motion of the electrons as they orbit the nucleus, c/1. Such a substance, subjected to a magnetic field, tends to align itself, c/2, so that a sheet of current circulates around the externally applied field. Figure c/3 is closely analogous to figure a/3; in the central gray area, the atomic currents cancel out, but the atoms at the outer surface form a sheet of bound current. However, whereas like charges repel and opposite charges attract, it works the other way around for currents: currents in the same direction attract, and currents in opposite directions repel. Therefore the bound currents in a material inserted inside a solenoid tend to *reinforce* the free currents, and the result is to strengthen the field. The total current is $I = I_{\text{free}} + I_{\text{bound}}$, and we define an altered version of the magnetic field,

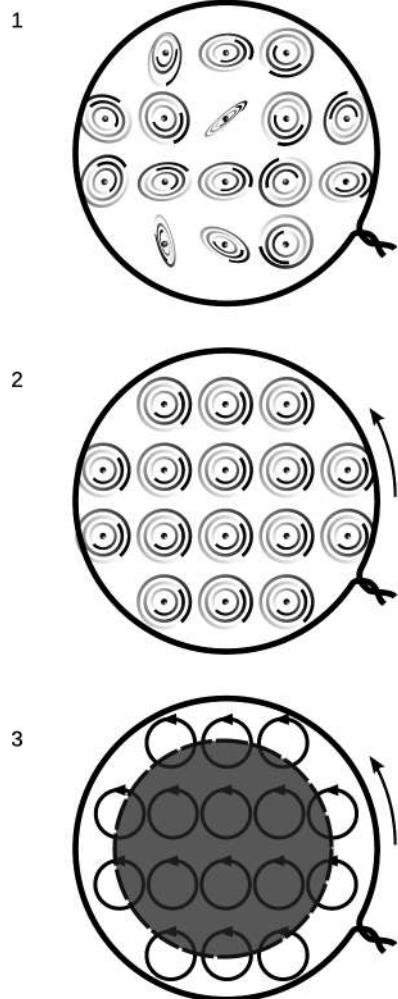
$$\mathbf{H} = \frac{\mathbf{B}}{\mu},$$

and rewrite Ampère's law as

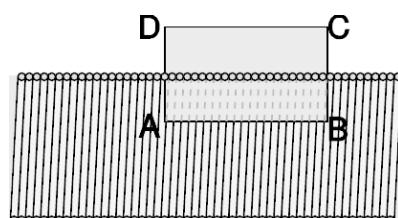
$$\Gamma_H = I_{\text{through, free.}}$$

The constant μ is the permeability, with a vacuum value of $\mu_0 = 4\pi k/c^2$. Here are the magnetic permeabilities of some substances:

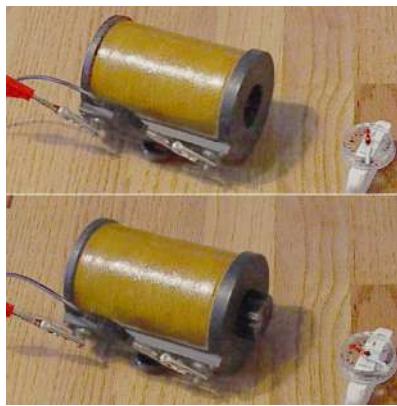
substance	μ/μ_0
vacuum	1
aluminum	1.00002
steel	700
transformer iron	4,000
mu-metal	20,000



c / The magnetic version of figure a. A magnetically permeable material is placed at the center of a solenoid.



d / Example 24: a cutaway view of a solenoid.



e / Example 24: without the iron core, the field is so weak that it barely deflects the compass. With it, the deflection is nearly 90° .



f / A transformer with a laminated iron core. The input and output coils are inside the paper wrapper. The iron core is the black part that passes through the coils at the center, and also wraps around them on the outside.



g / Example 25: ferrite beads. The top panel shows a clip-on type, while the bottom shows one built into a cable.

An iron-core electromagnet

example 24

▷ A solenoid has 1000 turns of wire wound along a cylindrical core with a length of 10 cm. If a current of 1.0 A is used, find the magnetic field inside the solenoid if the core is air, and if the core is made of iron with $\mu/\mu_0 = 4,000$.

▷ Air has essentially the same permability as vacuum, so using the result of example 13 on page 703, we find that the field is 0.013 T.

We now consider the case where the core is filled with iron. The original derivation in example 13 started from Ampère's law, which we now rewrite as $\Gamma_H = I_{\text{through, free}}$. As argued previously, the only significant contributions to the circulation come from line segment AB. This segment lies inside the iron, where $\mathbf{H} = \mathbf{B}/\mu$. The \mathbf{H} field is the same as in the air-core case, since the new form of Ampère's law only relates \mathbf{H} to the current in the wires (the free current). This means that $\mathbf{B} = \mu\mathbf{H}$ is greater by a factor of 4,000 than in the air-core case, or 52 T. This is an extremely intense field — so intense, in fact, that the iron's magnetic polarization would probably become saturated before we could actually get the field that high.

The electromagnet of example 24 could also be used as an inductor, and its inductance would be proportional to the permittivity of the core. This makes it possible to construct high-value inductors that are relatively compact. Permeable cores are also used in transformers.

A transformer or inductor with a permeable core does have some disadvantages, however, in certain applications. The oscillating magnetic field induces an electric field, and because the core is typically a metal, these currents dissipate energy strongly as heat. This behaves like a fairly large resistance in series with the coil. Figure f shows a method for reducing this effect. The iron core of this transformer has been constructed out of laminated layers, which has the effect of blocking the conduction of the eddy currents.

A ferrite bead

example 25

Cables designed to carry audio signals are typically made with two adjacent conductors, such that the current flowing out through one conductor comes back through the other one. Computer cables are similar, but usually have several such pairs bundled inside the insulator. This paired arrangement is known as differential mode, and has the advantage of cutting down on the reception and transmission of interference. In terms of transmission, the magnetic field created by the outgoing current is almost exactly canceled by the field from the return current, so electromagnetic waves are only weakly induced. In reception, both conductors are bathed in the same electric and magnetic fields, so an emf that adds current on one side subtracts current from the other side,

resulting in cancellation.

The opposite of differential mode is called common mode. In common mode, all conductors have currents flowing in the same direction. Even when a circuit is designed to operate in differential mode, it may not have exactly equal currents in the two conductors with $I_1 + I_2 = 0$, meaning that current is leaking off to ground at one end of the circuit or the other. Although paired cables are relatively immune to differential-mode interference, they do not have any automatic protection from common-mode interference.

Figure g shows a device for reducing common-mode interference called a ferrite bead, which surrounds the cable like a bead on a string. Ferrite is a magnetically permeable alloy. In this application, the ohmic properties of the ferrite actually turn out to be advantageous.

Let's consider common-mode transmission of interference. The bare cable has some DC resistance, but is also surrounded by a magnetic field, so it has inductance as well. This means that it behaves like a series L-R circuit, with an impedance that varies as $R + i\omega L$, where both R and L are very small. When we add the ferrite bead, the inductance is increased by orders of magnitude, but so is the resistance. Neither R nor L is actually constant with respect to frequency, but both are much greater than for the bare cable.

Suppose, for example, that a signal is being transmitted from a digital camera to a computer via a USB cable. The camera has an internal impedance that is on the order of 10Ω , the computer's input also has a $\sim 10 \Omega$ impedance, and in differential mode the ferrite bead has no effect, so the cable's impedance has its low, designed value (probably also about 10Ω , for good impedance matching). The signal is transmitted unattenuated from the camera to the computer, and there is almost no radiation from the cable.

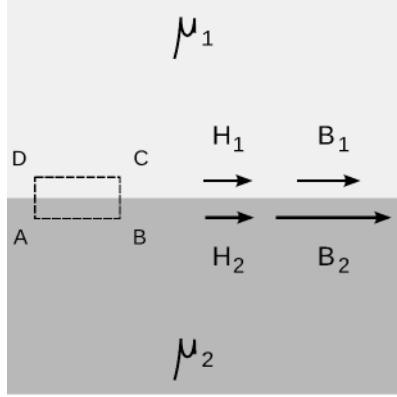
But in reality there will be a certain amount of common-mode current as well. With respect to common mode, the ferrite bead has a large impedance, with the exact value depending on frequency, but typically on the order of 100Ω for frequencies in the MHz range. We now have a series circuit consisting of three impedances: 10 , 100 , and 10Ω . For a given emf applied by an external radio wave, the current induced in the circuit has been attenuated by an order of magnitude, relative to its value without the ferrite bead.

Why is the ferrite necessary at all? Why not just insert ordinary air-core inductors in the circuit? We could, for example, have two solenoidal coils, one in the outgoing line and one in the return line, interwound with one another with their windings oriented so that

their differential-mode fields would cancel. There are two good reasons to prefer the ferrite bead design. One is that it allows a clip-on device like the one in the top panel of figure g, which can be added without breaking the circuit. The other is that our circuit will inevitably have some stray capacitance, and will therefore act like an LRC circuit, with a resonance at some frequency. At frequencies close to the resonant frequency, the circuit would absorb and transmit common-mode interference very strongly, which is exactly the opposite of the effect we were hoping to produce. The resonance peak could be made low and broad by adding resistance in series, but this extra resistance would attenuate the differential-mode signals as well as the common-mode ones. The ferrite's resistance, however, is actually a purely magnetic effect, so it vanishes in differential mode.

Surprisingly, some materials have magnetic permeabilities less than μ_0 . This cannot be accounted for in the model above, and although there are semiclassical arguments that can explain it to some extent, it is fundamentally a quantum mechanical effect. Materials with $\mu > \mu_0$ are called paramagnetic, while those with $\mu < \mu_0$ are referred to as diamagnetic. Diamagnetism is generally a much weaker effect than paramagnetism, and is easily masked if there is any trace of contamination from a paramagnetic material. Diamagnetic materials have the interesting property that they are repelled from regions of strong magnetic field, and it is therefore possible to levitate a diamagnetic object above a magnet, as in figure h.

A complete statement of Maxwell's equations in the presence of electric and magnetic materials is as follows:



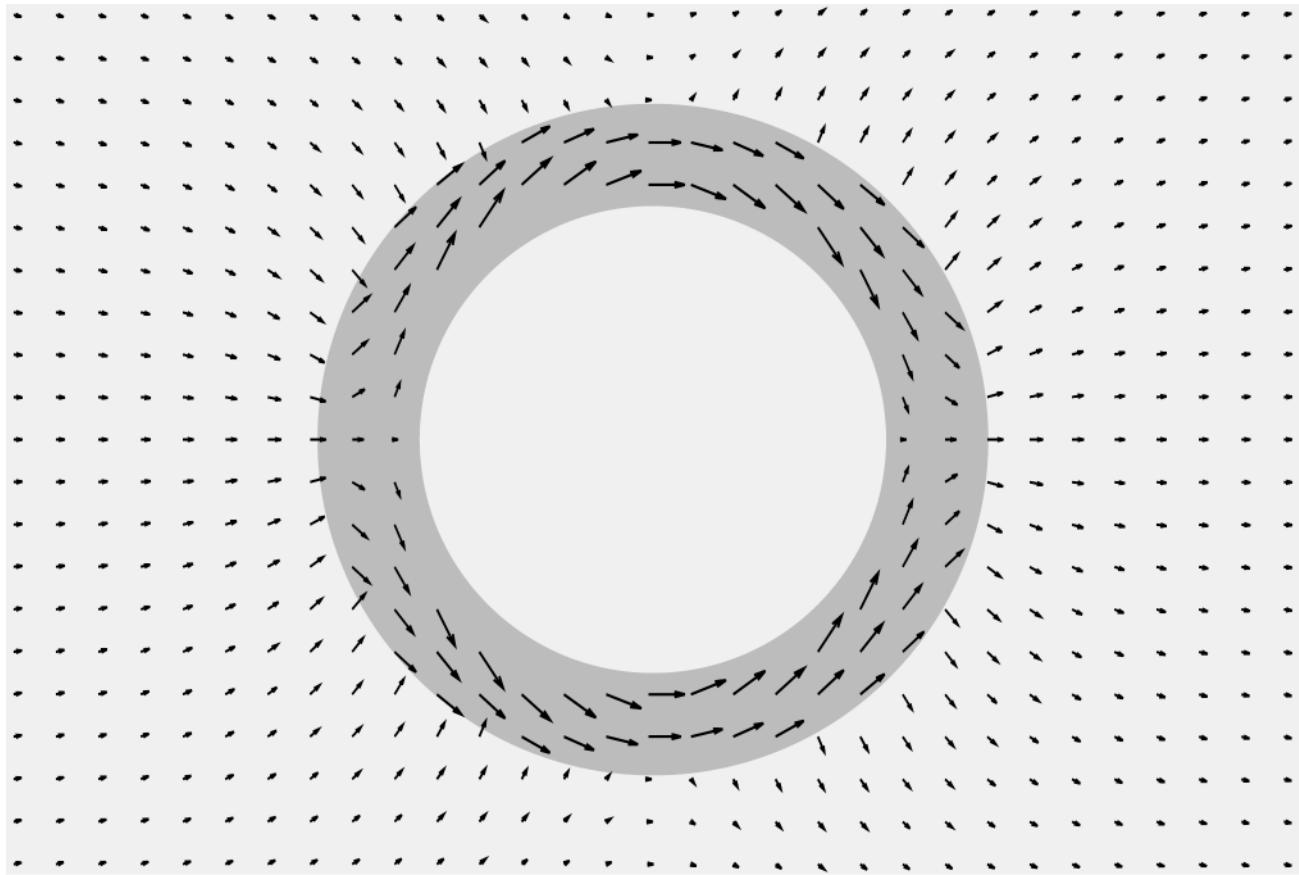
i / At a boundary between two substances with $\mu_2 > \mu_1$, the \mathbf{H} field has a continuous component parallel to the surface, which implies a discontinuity in the parallel component of the magnetic field \mathbf{B} .

$$\begin{aligned}\Phi_D &= q_{\text{free}} \\ \Phi_B &= 0 \\ \Gamma_E &= -\frac{d\Phi_B}{dt} \\ \Gamma_H &= \frac{d\Phi_D}{dt} + I_{\text{free}}\end{aligned}$$

Comparison with the vacuum case shows that the speed of an electromagnetic wave moving through a substance described by permittivity and permeability ϵ and μ is $1/\sqrt{\epsilon\mu}$. For most substances, $\mu \approx \mu_0$, and ϵ is highly frequency-dependent.

Suppose we have a boundary between two substances. By constructing a Gaussian or Ampèrean surface that extends across the boundary, we can arrive at various constraints on how the fields must behave as we move from one substance into the other, when there are no free currents or charges present, and the fields are static. An interesting example is the application of Faraday's law, $\Gamma_H = 0$, to the case where one medium — let's say it's air — has

a low permeability, while the other one has a very high one. We will violate Faraday's law unless the component of the \mathbf{H} field parallel to the boundary is a continuous function, $\mathbf{H}_{\parallel,1} = \mathbf{H}_{\parallel,2}$. This means that if μ/μ_0 is very high, the component of $\mathbf{B} = \mu\mathbf{H}$ parallel to the surface will have an abrupt discontinuity, being much stronger inside the high-permeability material. The result is that when a magnetic field enters a high-permeability material, it tends to twist abruptly to one side, and the pattern of the field tends to be channeled through the material like water through a funnel. In a transformer, a permeable core functions to channel more of the magnetic flux from the input coil to the output coil. Figure j shows another example, in which the effect is to shield the interior of the sphere from the externally imposed field. Special high-permeability alloys, with trade names like Mu-Metal, are sold for this purpose.



j / A hollow sphere with $\mu/\mu_0 = 10$, is immersed in a uniform, externally imposed magnetic field. The interior of the sphere is shielded from the field. The arrows map the magnetic field \mathbf{B} . (See homework problem 47, page 757.)

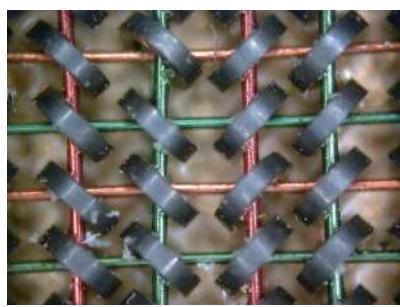
Ferromagnetism

The very last magnetic phenomenon we'll discuss is probably the very first experience you ever had of magnetism. Ferromagnetism is a phenomenon in which a material tends to organize itself so that it has a nonvanishing magnetic field. It is exhibited strongly by iron and nickel, which explains the origin of the name.

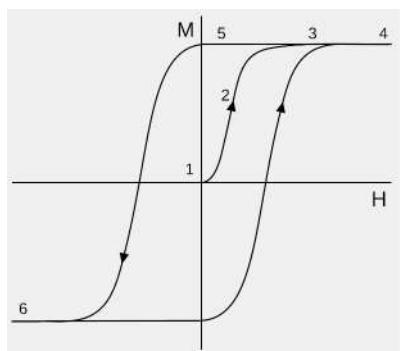
k / A model of ferromagnetism.



Figure k/1 is a simple one-dimensional model of ferromagnetism. Each magnetic compass needle represents an atom. The compasses in the chain are stable when aligned with one another, because each one's north end is attracted to its neighbor's south end. The chain can be turned around, k/2, without disrupting its organization, and the compasses do not realign themselves with the Earth's field, because their torques on one another are stronger than the Earth's torques on them. The system has a memory. For example, if I want to remind myself that my friend's address is 137 Coupling Ct., I can align the chain at an angle of 137 degrees. The model fails, however, as an explanation of real ferromagnetism, because in two or more dimensions, the most stable arrangement of a set of interacting magnetic dipoles is something more like k/3, in which alternating rows point in opposite directions. In this two-dimensional pattern, every compass is aligned in the most stable way with all four of its neighbors. This shows that ferromagnetism, like diamagnetism, has no purely classical explanation; a full explanation requires quantum mechanics.



l / Magnetic core memory.



m / A hysteresis curve.

Because ferromagnetic substances "remember" the history of how they were prepared, they are commonly used to store information in computers. Figure l shows 16 bits from an ancient (ca. 1970) 4-kilobyte random-access memory, in which each doughnut-shaped iron "core" can be magnetized in one of two possible directions, so that it stores one bit of information. Today, RAM is made of transistors rather than magnetic cores, but a remnant of the old technology remains in the term "core dump," meaning "memory dump," as in "my girlfriend gave me a total core dump about her mom's divorce." Most computer hard drives today do store their information on rotating magnetic platters, but the platter technology

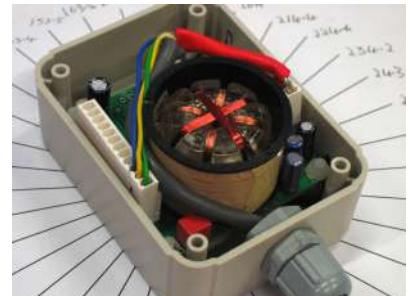
may be obsoleted by flash memory in the near future.

The memory property of ferromagnets can be depicted on the type of graph shown in figure m, known as a hysteresis curve. The y axis is the magnetization of a sample of the material — a measure of the extent to which its atomic dipoles are aligned with one another. If the sample is initially unmagnetized, 1, and a field H is externally applied, the magnetization increases, 2, but eventually becomes saturated, 3, so that higher fields do not result in any further magnetization, 4. The external field can then be reduced, 5, and even eliminated completely, but the material will retain its magnetization. It is a permanent magnet. To eliminate its magnetization completely, a substantial field must be applied in the opposite direction. If this reversed field is made stronger, then the substance will eventually become magnetized just as strongly in the opposite direction. Since the hysteresis curve is nonlinear, and is not a function (it has more than one value of M for a particular value of B), a ferromagnetic material does not have a single, well-defined value of the permeability μ ; a value like 4,000 for transformer iron represents some kind of a rough average.

The fluxgate compass

example 26

The fluxgate compass is a type of magnetic compass without moving parts, commonly used on ships and aircraft. An AC current is applied in a coil wound around a ferromagnetic core, driving the core repeatedly around a hysteresis loop. Because the hysteresis curve is highly nonlinear, the addition of an external field such as the Earth's alters the core's behavior. Suppose, for example, that the axis of the coil is aligned with the magnetic north-south. The core will reach saturation more quickly when the coil's field is in the same direction as the Earth's, but will not saturate as early in the next half-cycle, when the two fields are in opposite directions. With the use of multiple coils, the components of the Earth's field can be measured along two or three axes, permitting the compass's orientation to be determined in two or (for aircraft) three dimensions.



n / A fluxgate compass.

Sharp magnet poles

example 27

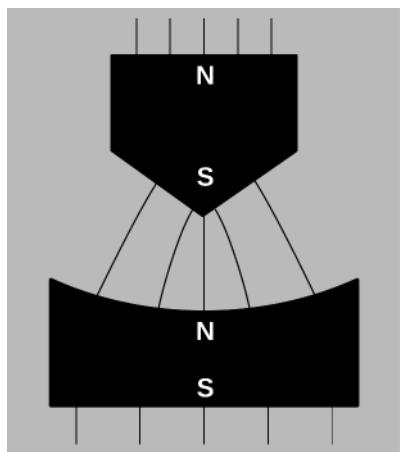
Although a ferromagnetic material does not really have a single value of the magnetic permeability, there is still a strong tendency to have $\mathbf{B}_{\parallel} \approx 0$ just outside the magnet's surface, for the same reasons as discussed above for high-permeability substances in general. For example, if we have a cylindrical bar magnet about the size and shape of your finger, magnetized lengthwise, then the field near the ends is nearly perpendicular to the surfaces, while the field near the sides, although it may be oriented nearly parallel to the surface, is very weak, so that we still have $\mathbf{B}_{\parallel} \approx 0$. This is in close analogy to the situation for the *electric* field near the surface of a conductor in equilibrium, for which $\mathbf{E}_{\parallel} = 0$.

This analogy is close enough so that we can recycle much of our knowledge about electrostatics.

For example, we saw in example 9, p. 546, and problem 37, p. 572, that charge tends to collect on the most highly curved portions of a conductor, and therefore becomes especially dense near a corner or knife-edge. This gives us a way of making especially intense magnetic fields. Most people would imagine that a very intense field could be made simply by using a very large and bulky permanent magnet, but this doesn't actually work very well, because magnetic dipole fields fall off as $1/r^3$, so that at a point near the surface, nearly all the field is contributed by atoms near the surface. Our analogy with electrostatics suggests that we should instead construct a permanent magnet with a sharp edge.

Figure o shows the cross-sectional shapes of two magnet poles used in the historic Stern-Gerlach experiment (sec. 14.1, p. 959). The external magnetic field is represented using field lines. The field lines enter and exit the surfaces perpendicularly, and they are particularly dense near the corner of the upper pole, indicating a strong field. The spreading of the field lines indicates that the field is strongly nonuniform, becoming much weaker toward the bottom of the gap between the poles. This strong nonuniformity was crucial for the experiment, in which the magnets were used as part of a dipole spectrometer. See example 7 and figure p on p. 591 for an explanation of an electric version of such a spectrometer.

This chapter is summarized on page 1089. Notation and terminology are tabulated on pages 1070-1071.



o / Example 27.

Problems

The symbols \checkmark , \blacksquare , etc. are explained on page 761.

1 A particle with a charge of 1.0 C and a mass of 1.0 kg is observed moving past point P with a velocity $(1.0 \text{ m/s})\hat{x}$. The electric field at point P is $(1.0 \text{ V/m})\hat{y}$, and the magnetic field is $(2.0 \text{ T})\hat{y}$. Find the force experienced by the particle. \checkmark \blacksquare

2 For a positively charged particle moving through a magnetic field, the directions of the \mathbf{v} , \mathbf{B} , and \mathbf{F} vectors are related by a right-hand rule:

- v** along the fingers, with the hand flat
- B** along the fingers, with the knuckles bent
- F** along the thumb

Make a three-dimensional model of the three vectors using pencils or rolled-up pieces of paper to represent the vectors assembled with their tails together. Make all three vectors perpendicular to each other. Now write down every possible way in which the rule could be rewritten by scrambling up the three symbols **v**, **B**, and **F**. Referring to your model, which are correct and which are incorrect? \blacksquare

3 A charged particle is released from rest. We see it start to move, and as it gets going, we notice that its path starts to curve. Can we tell whether this region of space has $\mathbf{E} \neq 0$, or $\mathbf{B} \neq 0$, or both? Assume that no other forces are present besides the possible electrical and magnetic ones, and that the fields, if they are present, are uniform. \blacksquare

4 A charged particle is in a region of space in which there is a uniform magnetic field $\mathbf{B} = B\hat{z}$. There is no electric field, and no other forces act on the particle. In each case, describe the future motion of the particle, given its initial velocity.

- (a) $\mathbf{v}_o = 0$
- (b) $\mathbf{v}_o = (1 \text{ m/s})\hat{z}$
- (c) $\mathbf{v}_o = (1 \text{ m/s})\hat{y}$ \blacksquare

5 (a) A line charge, with charge per unit length λ , moves at velocity v along its own length. How much charge passes a given point in time dt ? What is the resulting current?

\triangleright Answer, p. 1069

(b) Show that the units of your answer in part a work out correctly.

Remark: This constitutes a physical model of an electric current, and it would be a physically realistic model of a beam of particles moving in a vacuum, such as the electron beam in a television tube. It is not a physically realistic model of the motion of the electrons in a current-carrying wire, or of the ions in your nervous system; the motion of the charge carriers in these systems is much more complicated and chaotic, and there are charges of both signs, so that the total charge is zero. But even when the model is physically unrealistic, it still gives the right answers when you use it to compute magnetic effects. This is a remarkable fact, which we will not prove. The interested reader is referred to E.M. Purcell,

6 Two parallel wires of length L carry currents I_1 and I_2 . They are separated by a distance R , and we assume R is much less than L , so that our results for long, straight wires are accurate. The goal of this problem is to compute the magnetic forces acting between the wires.

(a) Neither wire can make a force on *itself*. Therefore, our first step in computing wire 1's force on wire 2 is to find the magnetic field made only by wire 1, in the space *occupied* by wire 2. Express this field in terms of the given quantities. ✓

(b) Let's model the current in wire 2 by pretending that there is a line charge inside it, possessing density per unit length λ_2 and moving at velocity v_2 . Relate λ_2 and v_2 to the current I_2 , using the result of problem 5a. Now find the magnetic force wire 1 makes on wire 2, in terms of I_1 , I_2 , L , and R . ▷ Answer, p. 1069

(c) Show that the units of the answer to part b work out to be newtons. ■

7 Suppose a charged particle is moving through a region of space in which there is an electric field perpendicular to its velocity vector, and also a magnetic field perpendicular to both the particle's velocity vector and the electric field. Show that there will be one particular velocity at which the particle can be moving that results in a total force of zero on it; this requires that you analyze both the magnitudes and the directions of the forces compared to one another. Relate this velocity to the magnitudes of the electric and magnetic fields. (Such an arrangement, called a velocity filter, is one way of determining the speed of an unknown particle.) ■

8 The following data give the results of two experiments in which charged particles were released from the same point in space, and the forces on them were measured:

$$\begin{aligned} q_1 &= 1 \text{ } \mu\text{C}, & \mathbf{v}_1 &= (1 \text{ m/s})\hat{\mathbf{x}}, & \mathbf{F}_1 &= (-1 \text{ mN})\hat{\mathbf{y}} \\ q_2 &= -2 \text{ } \mu\text{C}, & \mathbf{v}_2 &= (-1 \text{ m/s})\hat{\mathbf{x}}, & \mathbf{F}_2 &= (-2 \text{ mN})\hat{\mathbf{y}} \end{aligned}$$

The data are insufficient to determine the magnetic field vector; demonstrate this by giving two different magnetic field vectors, both of which are consistent with the data. ■

9 The following data give the results of two experiments in which charged particles were released from the same point in space, and the forces on them were measured:

$$\begin{aligned} q_1 &= 1 \text{ nC}, & \mathbf{v}_1 &= (1 \text{ m/s})\hat{\mathbf{z}}, & \mathbf{F}_1 &= (5 \text{ pN})\hat{\mathbf{x}} + (2 \text{ pN})\hat{\mathbf{y}} \\ q_2 &= 1 \text{ nC}, & \mathbf{v}_2 &= (3 \text{ m/s})\hat{\mathbf{z}}, & \mathbf{F}_2 &= (10 \text{ pN})\hat{\mathbf{x}} + (4 \text{ pN})\hat{\mathbf{y}} \end{aligned}$$

Is there a nonzero electric field at this point? A nonzero magnetic field? ■

10 This problem is a continuation of problem 6. Note that the answer to problem 6b is given on page 1069.

(a) Interchanging the 1's and 2's in the answer to problem 6b, what

is the magnitude of the magnetic force from wire 2 acting on wire 1? Is this consistent with Newton's third law?

(b) Suppose the currents are in the same direction. Make a sketch, and use the right-hand rule to determine whether wire 1 pulls wire 2 towards it, or pushes it away.

(c) Apply the right-hand rule again to find the direction of wire 2's force on wire 1. Does this agree with Newton's third law?

(d) What would happen if wire 1's current was in the opposite direction compared to wire 2's? ■

11 (a) In the photo of the vacuum tube apparatus in figure o on page 684, infer the direction of the magnetic field from the motion of the electron beam. (The answer is given in the answer to the self-check on that page.)

(b) Based on your answer to part a, find the direction of the currents in the coils.

(c) What direction are the electrons in the coils going?

(d) Are the currents in the coils repelling the currents consisting of the beam inside the tube, or attracting them? Check your answer by comparing with the result of problem 10. ■

12 A charged particle of mass m and charge q moves in a circle due to a uniform magnetic field of magnitude B , which points perpendicular to the plane of the circle.

(a) Assume the particle is positively charged. Make a sketch showing the direction of motion and the direction of the field, and show that the resulting force is in the right direction to produce circular motion.

(b) Find the radius, r , of the circle, in terms of m , q , v , and B . ✓

(c) Show that your result from part b has the right units.

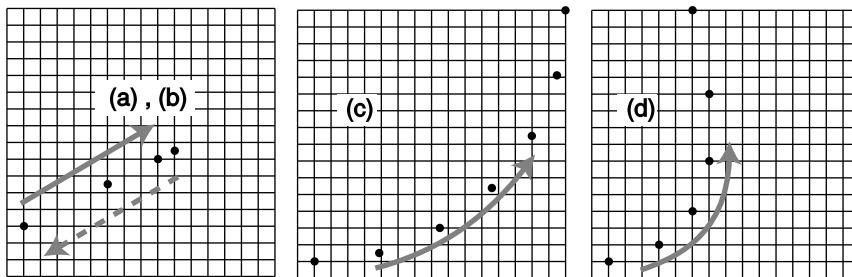
(d) Discuss all four variables occurring on the right-hand side of your answer from part b. Do they make sense? For instance, what should happen to the radius when the magnetic field is made stronger? Does your equation behave this way?

(e) Restate your result so that it gives the particle's angular frequency, ω , in terms of the other variables, and show that v drops out. ✓

Remark: A charged particle can be accelerated in a circular device called a cyclotron, in which a magnetic field is what keeps them from going off straight. This frequency is therefore known as the cyclotron frequency. The particles are accelerated by other forces (electric forces), which are AC. As long as the electric field is operated at the correct cyclotron frequency for the type of particles being manipulated, it will stay in sync with the particles, giving them a shove in the right direction each time they pass by. The particles are speeding up, so this only works because the cyclotron frequency is independent of velocity. ■

13 Each figure represents the motion of a positively charged particle. The dots give the particles' positions at equal time intervals. In each case, determine whether the motion was caused by an electric force, a magnetic force, or a frictional force, and explain

your reasoning. If possible, determine the direction of the magnetic or electric field. All fields are uniform. In (a), the particle stops for an instant at the upper right, but then comes back down and to the left, retracing the same dots. In (b), it stops on the upper right and stays there.



Problem 13.

14 One model of the hydrogen atom has the electron circling around the proton at a speed of 2.2×10^6 m/s, in an orbit with a radius of 0.05 nm. (Although the electron and proton really orbit around their common center of mass, the center of mass is very close to the proton, since it is 2000 times more massive. For this problem, assume the proton is stationary.) In homework problem 15, p. 567, you calculated the electric current created.

- (a) Now estimate the magnetic field created at the center of the atom by the electron. We are treating the circling electron as a current loop, even though it's only a single particle. ✓
- (b) Does the proton experience a nonzero force from the electron's magnetic field? Explain.
- (c) Does the electron experience a magnetic field from the proton? Explain.
- (d) Does the electron experience a magnetic field created by its own current? Explain.
- (e) Is there an electric force acting between the proton and electron? If so, calculate it. ✓
- (f) Is there a gravitational force acting between the proton and electron? If so, calculate it.
- (g) An inward force is required to keep the electron in its orbit – otherwise it would obey Newton's first law and go straight, leaving the atom. Based on your answers to the previous parts, which force or forces (electric, magnetic and gravitational) contributes significantly to this inward force?

[Based on a problem by Arnold Arons.]

15 The equation $B_z = \beta k I A / c^2 r^3$ was found on page 694 for the distant field of a dipole. Show, as asserted there, that the constant

β must be unitless. ■

16 The following data give the results of three experiments in which charged particles were released from the same point in space, and the forces on them were measured:

$$\begin{aligned}q_1 &= 1 \text{ C}, & \mathbf{v}_1 &= 0, & \mathbf{F}_1 &= (1 \text{ N})\hat{\mathbf{y}} \\q_2 &= 1 \text{ C}, & \mathbf{v}_2 &= (1 \text{ m/s})\hat{\mathbf{x}}, & \mathbf{F}_2 &= (1 \text{ N})\hat{\mathbf{y}} \\q_3 &= 1 \text{ C}, & \mathbf{v}_3 &= (1 \text{ m/s})\hat{\mathbf{z}}, & \mathbf{F}_3 &= 0\end{aligned}$$

Determine the electric and magnetic fields. ✓ ■

17 If you put four times more current through a solenoid, how many times more energy is stored in its magnetic field? ✓ ■

18 A Helmholtz coil is defined as a pair of identical circular coils lying in parallel planes and separated by a distance, h , equal to their radius, b . (Each coil may have more than one turn of wire.) Current circulates in the same direction in each coil, so the fields tend to reinforce each other in the interior region. This configuration has the advantage of being fairly open, so that other apparatus can be easily placed inside and subjected to the field while remaining visible from the outside. The choice of $h = b$ results in the most uniform possible field near the center. A photograph of a Helmholtz coil is shown in example 4 on page 684.

(a) Find the percentage drop in the field at the center of one coil, compared to the full strength at the center of the whole apparatus. ✓

(b) What value of h (not equal to b) would make this difference equal to zero? ✓ ■

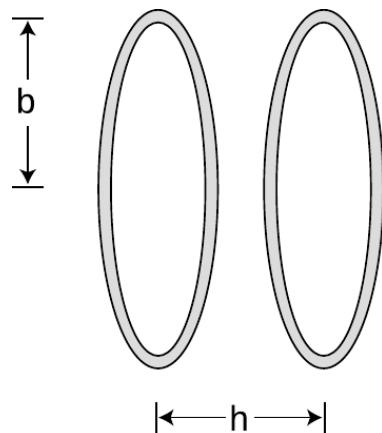
19 The figure shows a nested pair of circular wire loops used to create magnetic fields. (The twisting of the leads is a practical trick for reducing the magnetic fields they contribute, so the fields are very nearly what we would expect for an ideal circular current loop.) The coordinate system below is to make it easier to discuss directions in space. One loop is in the $y - z$ plane, the other in the $x - y$ plane. Each of the loops has a radius of 1.0 cm, and carries 1.0 A in the direction indicated by the arrow.

(a) Calculate the magnetic field that would be produced by *one* such loop, at its center. ✓

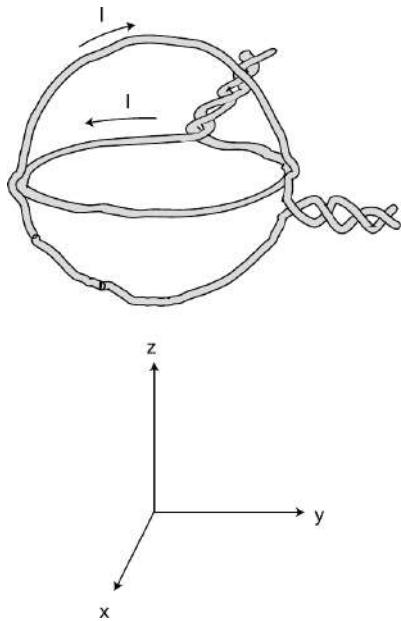
(b) Describe the direction of the magnetic field that would be produced, at its center, by the loop in the $x - y$ plane alone.

(c) Do the same for the other loop.

(d) Calculate the magnitude of the magnetic field produced by the two loops in combination, at their common center. Describe its direction. ✓ ■

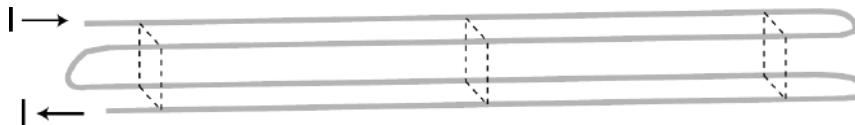


Problem 18.



Problem 19.

20 Four long wires are arranged, as shown, so that their cross-section forms a square, with connections at the ends so that current



Problem 20.

flows through all four before exiting. Note that the current is to the right in the two back wires, but to the left in the front wires. If the dimensions of the cross-sectional square (height and front-to-back) are b , find the magnetic field (magnitude and direction) along the long central axis. \checkmark ■

21 In problem 16, the three experiments gave enough information to determine both fields. Is it possible to design a procedure so that, using only two such experiments, we can always find \mathbf{E} and \mathbf{B} ? If so, design it. If not, why not? ■

22 Use the Biot-Savart law to derive the magnetic field of a long, straight wire, and show that this reproduces the result of example 6 on page 686. ■

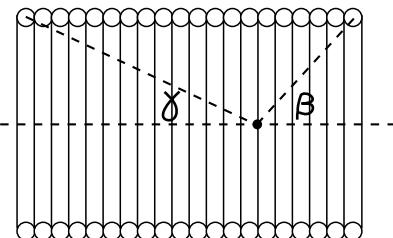
23 (a) Modify the calculation on page 691 to determine the component of the magnetic field of a sheet of charge that is perpendicular to the sheet. \checkmark

(b) Show that your answer has the right units.

(c) Show that your answer approaches zero as z approaches infinity.

(d) What happens to your answer in the case of $a = b$? Explain why this makes sense. ■

24 Consider two solenoids, one of which is smaller so that it can be put inside the other. Assume they are long enough so that each one only contributes significantly to the field inside itself, and the interior fields are nearly uniform. Consider the configuration where the small one is inside the big one with their currents circulating in the same direction, and a second configuration in which the currents circulate in opposite directions. Compare the energies of these configurations with the energy when the solenoids are far apart. Based on this reasoning, which configuration is stable, and in which configuration will the little solenoid tend to get twisted around or spit out? \triangleright Hint, p. 1037 ■



Problem 25.

25 (a) A solenoid can be imagined as a series of circular current loops that are spaced along their common axis. Integrate the result of example 12 on page 700 to show that the field on the axis of a solenoid can be written as $B = (2\pi k\eta/c^2)(\cos \beta + \cos \gamma)$, where the angles β and γ are defined in the figure.

(b) Show that in the limit where the solenoid is very long, this exact result agrees with the approximate one derived in example 13 on page 703 using Ampère's law.

- (c) Note that, unlike the calculation using Ampère's law, this one is valid at points that are near the mouths of the solenoid, or even outside it entirely. If the solenoid is long, at what point on the axis is the field equal to one half of its value at the center of the solenoid?
 (d) What happens to your result when you apply it to points that are very far away from the solenoid? Does this make sense? ■

26 The first step in the proof of Ampère's law on page 704 is to show that Ampère's law holds in the case shown in figure f/1, where a circular Ampèrian loop is centered on a long, straight wire that is perpendicular to the plane of the loop. Carry out this calculation, using the result for the field of a wire that was established without using Ampère's law. ■

27 A certain region of space has a magnetic field given by $\mathbf{B} = bx\hat{\mathbf{y}}$. Find the electric current flowing through the square defined by $z = 0$, $0 \leq x \leq a$, and $0 \leq y \leq a$. ✓ ■

28 Perform a calculation similar to the one in problem 54, but for a logarithmic spiral, defined by $r = we^{u\theta}$, and show that the field is $B = (kI/c^2u)(1/a - 1/b)$. Note that the solution to problem 54 is given in the back of the book. ■

29 (a) For the geometry described in example 8 on page 689, find the field at a point the lies in the plane of the wires, but not between the wires, at a distance b from the center line. Use the same technique as in that example.

(b) Now redo the calculation using the technique demonstrated on page 694. The integrals are nearly the same, but now the reasoning is reversed: you already know $\beta = 1$, and you want to find an unknown field. The only difference in the integrals is that you are tiling a different region of the plane in order to mock up the currents in the two wires. Note that you can't tile a region that contains a point of interest, since the technique uses the field of a distant dipole. ✓ ■

30 (a) A long, skinny solenoid consists of N turns of wire wrapped uniformly around a hollow cylinder of length ℓ and cross-sectional area A . Find its inductance. ✓

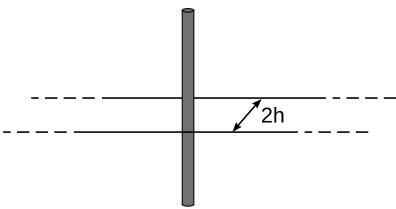
(b) Show that your answer has the right units to be an inductance. ■

31 Consider two solenoids, one of which is smaller so that it can be put inside the other. Assume they are long enough to act like ideal solenoids, so that each one only contributes significantly to the field inside itself, and the interior fields are nearly uniform. Consider the configuration where the small one is partly inside and partly hanging out of the big one, with their currents circulating in the same direction. Their axes are constrained to coincide.

(a) Find the difference in the magnetic energy between the configuration where the solenoids are separate and the configuration where



A nautilus shell is approximately a logarithmic spiral, of the type in problem 28.



Problem 32.

the small one is inserted into the big one. Your equation will include the length x of the part of the small solenoid that is inside the big one, as well as other relevant variables describing the two solenoids. ✓

(b) Based on your answer to part a, find the force acting ■

32 Verify Ampère's law in the case shown in the figure, assuming the known equation for the field of a wire. A wire carrying current I passes perpendicularly through the center of the rectangular Ampèrian surface. The length of the rectangle is infinite, so it's not necessary to compute the contributions of the ends. ■

33 The purpose of this problem is to find how the gain of a transformer depends on its construction.

(a) The number of loops of wire, N , in a solenoid is changed, while keeping the length constant. How does the impedance depend on N ? State your answer as a proportionality, e.g., $Z \propto N^3$ or $Z \propto N^{-5}$.

(b) For a given AC voltage applied across the inductor, how does the magnetic field depend on N ? You need to take into account both the dependence of a solenoid's field on N for a given current and your answer to part a, which affects the current.

(c) Now consider a transformer consisting of two solenoids. The input side has N_1 loops, and the output N_2 . We wish to find how the output voltage V_2 depends on N_1 , N_2 , and the input voltage V_1 . The text has already established $V_2 \propto V_1 N_2$, so it only remains to find the dependence on N_1 . Use your result from part b to accomplish this. The ratio V_2/V_1 is called the voltage gain. ■

34 Problem 33 dealt with the dependence of a transformer's gain on the number of loops of wire in the input solenoid. Carry out a similar analysis of how the gain depends on the frequency at which the circuit is operated. ■

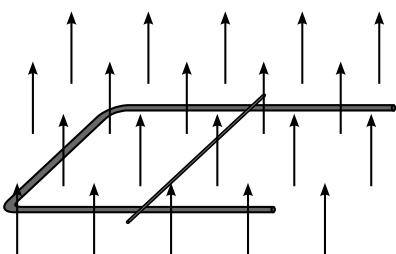
35 A U-shaped wire makes electrical contact with a second, straight wire, which rolls along it to the right, as shown in the figure. The whole thing is immersed in a uniform magnetic field, which is perpendicular to the plane of the circuit. The resistance of the rolling wire is much greater than that of the U.

(a) Find the direction of the force on the wire based on conservation of energy.

(b) Verify the direction of the force using right-hand rules.

(c) Find the magnitude of the force acting on the wire. There is more than one way to do this, but please do it using Faraday's law (which works even though it's the Ampèrian surface itself that is changing, rather than the field). ✓

(d) Consider how the answer to part a would have changed if the direction of the field had been reversed, and also do the case where the direction of the rolling wire's motion is reversed. Verify that this is in agreement with your answer to part c. ■

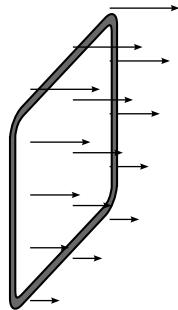


Problem 35.

36 A charged particle is in motion at speed v , in a region of vacuum through which an electromagnetic wave is passing. In what direction should the particle be moving in order to minimize the total force acting on it? Consider both possibilities for the sign of the charge. (Based on a problem by David J. Raymond.) ■

37 A wire loop of resistance R and area A , lying in the $y - z$ plane, falls through a nonuniform magnetic field $\mathbf{B} = kz\hat{\mathbf{x}}$, where k is a constant. The z axis is vertical.

- (a) Find the direction of the force on the wire based on conservation of energy.
- (b) Verify the direction of the force using right-hand rules.
- (c) Find the magnetic force on the wire. ✓ ■



Problem 37.

38 A capacitor has parallel plates of area A , separated by a distance h . If there is a vacuum between the plates, then Gauss's law gives $E = 4\pi k\sigma = 4\pi kq/A$ for the field between the plates, and combining this with $E = V/h$, we find $C = q/V = (1/4\pi k)A/h$.

- (a) Generalize this derivation to the case where there is a dielectric between the plates.
- (b) Suppose we have a list of possible materials we could choose as dielectrics, and we wish to construct a capacitor that will have the highest possible energy density, U_e/v , where v is the volume. For each dielectric, we know its permittivity ϵ , and also the maximum electric field E it can sustain without breaking down and allowing sparks to cross between the plates. Write the maximum energy density in terms of these two variables, and determine a figure of merit that could be used to decide which material would be the best choice. ■

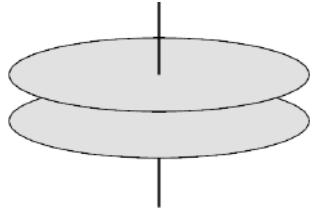
39 (a) For each term appearing on the right side of Maxwell's equations, give an example of an everyday situation it describes.

- (b) Most people doing calculations in the SI system of units don't use k and k/c^2 . Instead, they express everything in terms of the constants

$$\epsilon_0 = \frac{1}{4\pi k} \quad \text{and}$$

$$\mu_0 = \frac{4\pi k}{c^2}.$$

Rewrite Maxwell's equations in terms of these constants, eliminating k and c everywhere. ■



Problem 42.

40 (a) Prove that in an electromagnetic plane wave, half the energy is in the electric field and half in the magnetic field.

(b) Based on your result from part a, find the proportionality constant in the relation $d\mathbf{p} \propto \mathbf{E} \times \mathbf{B} dv$, where $d\mathbf{p}$ is the momentum of the part of a plane light wave contained in the volume dv . The vector $\mathbf{E} \times \mathbf{B}$, multiplied by the appropriate constant, is known as the Poynting vector, and even outside the context of an electromagnetic plane wave, it can be interpreted as a momentum density or rate of energy flow. (To do this problem, you need to know the relativistic relationship between the energy and momentum of a beam of light from problem 11 on p. 460.) \checkmark ■

41 (a) A beam of light has cross-sectional area A and power P , i.e., P is the number of joules per second that enter a window through which the beam passes. Find the energy density U/v in terms of P , A , and universal constants.

(b) Find $\tilde{\mathbf{E}}$ and $\tilde{\mathbf{B}}$, the amplitudes of the electric and magnetic fields, in terms of P , A , and universal constants (i.e., your answer should *not* include U or v). You will need the result of problem 40a. A real beam of light usually consists of many short wavetrains, not one big sine wave, but don't worry about that. \checkmark ▷ Hint, p. 1037

(c) A beam of sunlight has an intensity of $P/A = 1.35 \times 10^3 \text{ W/m}^2$, assuming no clouds or atmospheric absorption. This is known as the solar constant. Compute $\tilde{\mathbf{E}}$ and $\tilde{\mathbf{B}}$, and compare with the strengths of static fields you experience in everyday life: $E \sim 10^6 \text{ V/m}$ in a thunderstorm, and $B \sim 10^{-3} \text{ T}$ for the Earth's magnetic field. \checkmark ■

42 The circular parallel-plate capacitor shown in the figure is being charged up over time, with the voltage difference across the plates varying as $V = st$, where s is a constant. The plates have radius b , and the distance between them is d . We assume $d \ll b$, so that the electric field between the plates is uniform, and parallel to the axis. Find the induced magnetic field at a point between the plates, at a distance R from the axis. ▷ Hint, p. 1037 \checkmark ■

43 A positively charged particle is released from rest at the origin at $t = 0$, in a region of vacuum through which an electromagnetic wave is passing. The particle accelerates in response to the wave. In this region of space, the wave varies as $\mathbf{E} = \hat{\mathbf{x}}\tilde{E} \sin \omega t$, $\mathbf{B} = \hat{\mathbf{y}}\tilde{B} \sin \omega t$, and we assume that the particle has a relatively large value of m/q , so that its response to the wave is sluggish, and it never ends up moving at any speed comparable to the speed of light. Therefore we don't have to worry about the spatial variation of the wave; we can just imagine that these are uniform fields imposed by some external mechanism on this region of space.

(a) Find the particle's coordinates as functions of time. \checkmark

(b) Show that the motion is confined to $-z_{max} \leq z \leq z_{max}$, where $z_{max} = 1.101 \left(q^2 \tilde{E} \tilde{B} / m^2 \omega^3 \right)$. ■

44 Electromagnetic waves are supposed to have their electric and magnetic fields perpendicular to each other. (Throughout this problem, assume we're talking about waves traveling through a vacuum, and that there is only a single sine wave traveling in a single direction, not a superposition of sine waves passing through each other.) Suppose someone claims they can make an electromagnetic wave in which the electric and magnetic fields lie in the same plane. Prove that this is impossible based on Maxwell's equations. ■

45 Repeat the self-check on page 739, but with one change in the procedure: after we charge the capacitor, we open the circuit, and then continue with the observations. ■

46 On page 742, I proved that $\mathbf{H}_{\parallel,1} = \mathbf{H}_{\parallel,2}$ at the boundary between two substances if there is no free current and the fields are static. In fact, each of Maxwell's four equations implies a constraint with a similar structure. Some are constraints on the field components parallel to the boundary, while others are constraints on the perpendicular parts. Since some of the fields referred to in Maxwell's equations are the electric and magnetic fields \mathbf{E} and \mathbf{B} , while others are the auxiliary fields \mathbf{D} and \mathbf{H} , some of the constraints deal with \mathbf{E} and \mathbf{B} , others with \mathbf{D} and \mathbf{H} . Find the other three constraints. ■

47 (a) Figure j on page 743 shows a hollow sphere with $\mu/\mu_0 = x$, inner radius a , and outer radius b , which has been subjected to an external field \mathbf{B}_o . Finding the fields on the exterior, in the shell, and on the interior requires finding a set of fields that satisfies five boundary conditions: (1) far from the sphere, the field must approach the constant \mathbf{B}_o ; (2) at the outer surface of the sphere, the field must have $\mathbf{H}_{\parallel,1} = \mathbf{H}_{\parallel,2}$, as discussed on page 742; (3) the same constraint applies at the inner surface of the sphere; (4) and (5) there is an additional constraint on the fields at the inner and outer surfaces, as found in problem 46. The goal of this problem is to find the solution for the fields, and from it, to prove that the interior field is uniform, and given by

$$\mathbf{B} = \left[\frac{9x}{(2x+1)(x+2) - 2\frac{a^3}{b^3}(x-1)^2} \right] \mathbf{B}_o.$$

This is a very difficult problem to solve from first principles, because it's not obvious what form the fields should have, and if you hadn't been told, you probably wouldn't have guessed that the interior field would be uniform. We could, however, guess that once the sphere becomes polarized by the external field, it would become a dipole, and at $r \gg b$, the field would be a uniform field superimposed on the field of a dipole. It turns out that even close to the sphere, the solution has exactly this form. In order to complete the solution, we need to find the field in the shell ($a < r < b$), but the only way this field could match up with the detailed angular variation of the

interior and exterior fields would be if it was also a superposition of a uniform field with a dipole field. The final result is that we have four unknowns: the strength of the dipole component of the external field, the strength of the uniform and dipole components of the field within the shell, and the strength of the uniform interior field. These four unknowns are to be determined by imposing constraints (2) through (5) above.

(b) Show that the expression from part a has physically reasonable behavior in its dependence on x and a/b . ■

48 Two long, parallel strips of thin metal foil form a configuration like a long, narrow sandwich. The air gap between them has height h , the width of each strip is w , and their length is ℓ . Each strip carries current I , and we assume for concreteness that the currents are in opposite directions, so that the magnetic force, F , between the strips is repulsive.

- (a) Find the force in the limit of $w \gg h$. ✓
- (b) Find the force in the limit of $w \ll h$, which is like two ordinary wires.
- (c) Discuss the relationship between the two results. ■

49 Suppose we are given a permanent magnet with a complicated, asymmetric shape. Describe how a series of measurements with a magnetic compass could be used to determine the strength and direction of its magnetic field at some point of interest. Assume that you are only able to see the direction to which the compass needle settles; you cannot measure the torque acting on it. ■

50 On page 709, the curl of $x\hat{y}$ was computed. Now consider the fields $x\hat{x}$ and $y\hat{y}$.

- (a) Sketch these fields.
- (b) Using the same technique of explicitly constructing a small square, prove that their curls are both zero. Do not use the component form of the curl; this was one step in *deriving* the component form of the curl. ■

51 If you watch a movie played backwards, some vectors reverse their direction. For instance, people walk backwards, with their velocity vectors flipped around. Other vectors, such as forces, keep the same direction, e.g., gravity still pulls down. An electric field is another example of a vector that doesn't turn around: positive charges are still positive in the time-reversed universe, so they still make diverging electric fields, and likewise for the converging fields around negative charges.

- (a) How does the momentum of a material object behave under time-reversal? ▷ Solution, p. 1049
- (b) The laws of physics are still valid in the time-reversed universe. For example, show that if two material objects are interacting, and momentum is conserved, then momentum is still conserved in the time-reversed universe. ▷ Solution, p. 1049

- (c) Discuss how currents and magnetic fields would behave under time reversal. \triangleright Hint, p. 1037
(d) Similarly, show that the equation $d\mathbf{p} \propto \mathbf{E} \times \mathbf{B}$ is still valid under time reversal. \blacksquare

52 This problem is a more advanced exploration of the time-reversal ideas introduced in problem 51.

- (a) In that problem, we assumed that charge did not flip its sign under time reversal. Suppose we make the opposite assumption, that charge *does* change its sign. This is an idea introduced by Richard Feynman: that antimatter is really matter traveling backward in time! Determine the time-reversal properties of \mathbf{E} and \mathbf{B} under this new assumption, and show that $d\mathbf{p} \propto \mathbf{E} \times \mathbf{B}$ is still valid under time-reversal.
(b) Show that Maxwell's equations are time-reversal symmetric, i.e., that if the fields $\mathbf{E}(x, y, z, t)$ and $\mathbf{B}(x, y, z, t)$ satisfy Maxwell's equations, then so do $\mathbf{E}(x, y, z, -t)$ and $\mathbf{B}(x, y, z, -t)$. Demonstrate this under both possible assumptions about charge, $q \rightarrow q$ and $q \rightarrow -q$. \blacksquare

53 The purpose of this problem is to prove that the constant of proportionality a in the equation $dU_m = aB^2 dv$, for the energy density of the magnetic field, is given by $a = c^2/8\pi k$ as asserted on page 693. The geometry we'll use consists of two sheets of current, like a sandwich with nothing in between but some vacuum in which there is a magnetic field. The currents are in opposite directions, and we can imagine them as being joined together at the ends to form a complete circuit, like a tube made of paper that has been squashed almost flat. The sheets have lengths L in the direction parallel to the current, and widths w . They are separated by a distance d , which, for convenience, we assume is small compared to L and w . Thus each sheet's contribution to the field is uniform, and can be approximated by the expression $2\pi k\eta/c^2$.

- (a) Make a drawing similar to the one in figure 11.2.1 on page 692, and show that in this opposite-current configuration, the magnetic fields of the two sheets reinforce in the region between them, producing double the field, but cancel on the outside.
(b) By analogy with the case of a single strand of wire, one sheet's force on the other is ILB_1 , where $I = \eta w$ is the total current in one sheet, and $B_1 = B/2$ is the field contributed by only one of the sheets, since the sheet can't make any net force on itself. Based on your drawing and the right-hand rule, show that this force is repulsive.

For the rest of the problem, consider a process in which the sheets start out touching, and are then separated to a distance d . Since the force between the sheets is repulsive, they do mechanical work on the outside world as they are separated, in much the same way that the piston in an engine does work as the gases inside the cylinder expand. At the same time, however, there is an induced emf

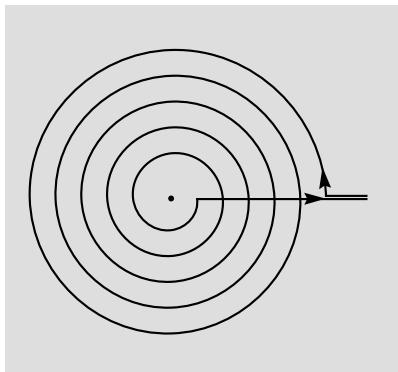
which would tend to extinguish the current, so in order to maintain a constant current, energy will have to be drained from a battery. There are three types of energy involved: the increase in the magnetic field energy, the increase in the energy of the outside world, and the decrease in energy as the battery is drained. (We assume the sheets have very little resistance, so there is no ohmic heating involved.) ✓

(c) Find the mechanical work done by the sheets, which equals the increase in the energy of the outside world. Show that your result can be stated in terms of η , the final volume $v = wLd$, and nothing else but numerical and physical constants. ✓

(d) The power supplied by the battery is $P = I\Gamma_E$ (like $P = I\Delta V$, but with an emf instead of a voltage difference), and the circulation is given by $\Gamma = -d\Phi_B/dt$. The negative sign indicates that the battery is being drained. Calculate the energy supplied by the battery, and, as in part c, show that the result can be stated in terms of η , v , and universal constants. ✓

(e) Find the increase in the magnetic-field energy, in terms of η , v , and the unknown constant a . ✓

(f) Use conservation of energy to relate your answers from parts c, d, and e, and solve for a . ✓ ■



Problem 54.

54 Magnet coils are often wrapped in multiple layers. The figure shows the special case where the layers are all confined to a single plane, forming a spiral. Since the thickness of the wires (plus their insulation) is fixed, the spiral that results is a mathematical type known as an Archimedean spiral, in which the turns are evenly spaced. The equation of the spiral is $r = w\theta$, where w is a constant. For a spiral that starts from $r = a$ and ends at $r = b$, show that the field at the center is given by $(kI/c^2w) \ln b/a$.

▷ Solution, p. 1049 ■

55 Resolve the following paradox. A capacitance C is initially charged, and is then connected to another capacitance C , forming a loop. With the charge now shared equally, the energy is halved. If the connection is made using wires that have finite resistance, then this energy loss could be explained through resistive heating. But how is conservation of energy satisfied if the resistance of the wires is zero? ■

56 The figure shows a simplified example of a device called a sector mass spectrometer. In an oven near the bottom, positively ionized atoms are produced. For simplicity, we assume that the atoms are all singly ionized. They may have different masses, however, and the goal is to separate them according to these masses. In the example shown in the figure, there are two different masses present. The reason this is called a “sector” mass spectrometer is that it contains two regions of uniform fields.

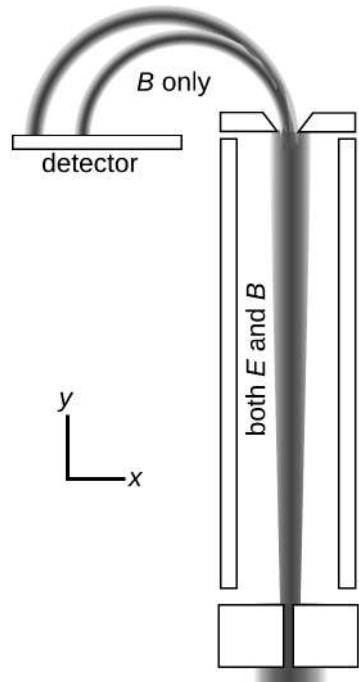
In the first sector, between the two long capacitor plates, there is an electric field E in the x direction. Superimposed on this is a uniform magnetic field B in the negative z direction (into the page). As analyzed in problem 7, these fields are chosen so that ions at a certain velocity v are not deflected. You will need the result of that problem in order to do this problem. Only the ions with the correct velocity make it out through the slits at the upper end of the capacitor.

In the second sector, at the top, there is no electric field, only a magnetic field, which we assume for simplicity to have the same magnitude and direction as in the first sector. This causes the beam to bend into a semicircular arc and hit a detector. In the first such spectrometers, this detector was simply some photographic film, whereas in modern ones it would probably be a silicon chip similar to the sensor of a camera.

The diameter h of the semicircle depends on the mass m of the ion. The quantity $\Delta h/\Delta m$ tells us how good the spectrometer is at separating similar masses.

- (a) Express $\Delta h/\Delta m$ in terms of E , B , and e , eliminating v (which we can neither control nor measure directly). ✓
- (b) Show that the units of your answer make sense.
- (c) You will have found that increasing E makes the spectrometer more sensitive, while increasing B makes it less so. Explain physically why this is so. What stops us from getting an arbitrarily large sensitivity simply by making B small enough?

Remark: This design makes inefficient use of the ion source’s intensity, because any ions with the wrong velocity are wasted. For this reason, real-world spectrometers of this type include complicated focusing elements. ■



Problem 56.

Key to symbols:

■ easy ■ typical ■ challenging ■ difficult ■ very difficult
✓ An answer check is available at www.lightandmatter.com.

Exercises

Exercise 11B: Polarization

Apparatus:

calcite (Iceland spar) crystal

polaroid film

1. Lay the crystal on a piece of paper that has print on it. You will observe a double image. See what happens if you rotate the crystal.

Evidently the crystal does something to the light that passes through it on the way from the page to your eye. One beam of light enters the crystal from underneath, but two emerge from the top; by conservation of energy the energy of the original beam must be shared between them. Consider the following three possible interpretations of what you have observed:

- (a) The two new beams differ from each other, and from the original beam, only in energy. Their other properties are the same.
- (b) The crystal adds to the light some mysterious new property (not energy), which comes in two flavors, X and Y. Ordinary light doesn't have any of either. One beam that emerges from the crystal has some X added to it, and the other beam has Y.
- (c) There is some mysterious new property that is possessed by all light. It comes in two flavors, X and Y, and most ordinary light sources make an equal mixture of type X and type Y light. The original beam is an even mixture of both types, and this mixture is then split up by the crystal into the two purified forms.

In parts 2 and 3 you'll make observations that will allow you to figure out which of these is correct.

2. Now place a polaroid film over the crystal and see what you observe. What happens when you rotate the film in the horizontal plane? Does this observation allow you to rule out any of the three interpretations?

3. Now put the polaroid film under the crystal and try the same thing. Putting together all your observations, which interpretation do you think is correct?

4. Look at an overhead light fixture through the polaroid, and try rotating it. What do you observe? What does this tell you about the light emitted by the lightbulb?

5. Now position yourself with your head under a light fixture and directly over a shiny surface, such as a glossy tabletop. You'll see the lamp's reflection, and the light coming from the lamp to your eye will have undergone a reflection through roughly a 180-degree angle (i.e., it very nearly reversed its direction). Observe this reflection through the polaroid, and try rotating it. Finally, position yourself so that you are seeing glancing reflections, and try the same thing. Summarize what happens to light with properties X and Y when it is reflected. (This is the principle behind polarizing sunglasses.)



Chapter 12

Optics

12.1 The ray model of light

Ads for one Macintosh computer bragged that it could do an arithmetic calculation in less time than it took for the light to get from the screen to your eye. We find this impressive because of the contrast between the speed of light and the speeds at which we interact with physical objects in our environment. Perhaps it shouldn't surprise us, then, that Newton succeeded so well in explaining the motion of objects, but was far less successful with the study of light.

The climax of our study of electricity and magnetism was discovery that light is an electromagnetic wave. Knowing this, however, is not the same as knowing everything about eyes and telescopes. In fact, the full description of light as a wave can be rather cumbersome. We will instead spend most of our treatment of optics making use of a simpler model of light, the ray model, which does a fine job in most practical situations. Not only that, but we will even backtrack a little and start with a discussion of basic ideas about light and vision that predated the discovery of electromagnetic waves.

12.1.1 The nature of light

The cause and effect relationship in vision

Despite its title, this chapter is far from your first look at light. That familiarity might seem like an advantage, but most people have never thought carefully about light and vision. Even smart people who have thought hard about vision have come up with incorrect ideas. The ancient Greeks, Arabs and Chinese had theories of light and vision, all of which were mostly wrong, and all of which were accepted for thousands of years.

One thing the ancients did get right is that there is a distinction between objects that emit light and objects that don't. When you see a leaf in the forest, it's because three different objects are doing their jobs: the leaf, the eye, and the sun. But luminous objects like the sun, a flame, or the filament of a light bulb can be seen by the eye without the presence of a third object. Emission of light is often, but not always, associated with heat. In modern times, we are familiar with a variety of objects that glow without being heated, including fluorescent lights and glow-in-the-dark toys.

How do we see luminous objects? The Greek philosophers Pythagoras (b. ca. 560 BC) and Empedocles of Acragas (b. ca. 492 BC), who unfortunately were very influential, claimed that when you looked at a candle flame, the flame and your eye were both sending out some kind of mysterious stuff, and when your eye's stuff collided with the candle's stuff, the candle would become evident to your sense of sight.

Bizarre as the Greek “collision of stuff theory” might seem, it had a couple of good features. It explained why both the candle and your eye had to be present for your sense of sight to function. The theory could also easily be expanded to explain how we see nonluminous objects. If a leaf, for instance, happened to be present at the site of the collision between your eye's stuff and the candle's stuff, then the leaf would be stimulated to express its green nature, allowing you to perceive it as green.

Modern people might feel uneasy about this theory, since it suggests that greenness exists only for our seeing convenience, implying a human precedence over natural phenomena. Nowadays, people would expect the cause and effect relationship in vision to be the other way around, with the leaf doing something to our eye rather than our eye doing something to the leaf. But how can you tell? The most common way of distinguishing cause from effect is to determine which happened first, but the process of seeing seems to occur too quickly to determine the order in which things happened. Certainly there is no obvious time lag between the moment when you move your head and the moment when your reflection in the mirror moves.

Today, photography provides the simplest experimental evidence that nothing has to be emitted from your eye and hit the leaf in order to make it “greenify.” A camera can take a picture of a leaf even if there are no eyes anywhere nearby. Since the leaf appears green regardless of whether it is being sensed by a camera, your eye, or an insect’s eye, it seems to make more sense to say that the leaf’s greenness is the cause, and something happening in the camera or eye is the effect.

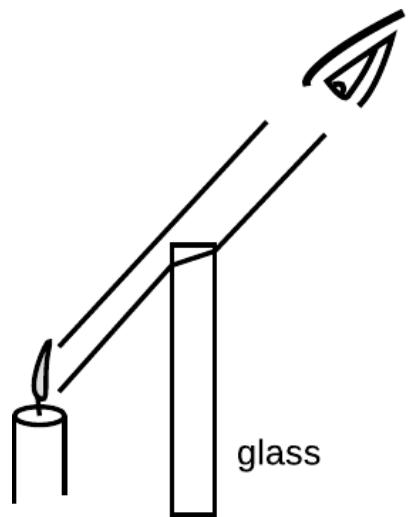
Light is a thing, and it travels from one point to another.

Another issue that few people have considered is whether a candle’s flame simply affects your eye directly, or whether it sends out light which then gets into your eye. Again, the rapidity of the effect makes it difficult to tell what’s happening. If someone throws a rock at you, you can see the rock on its way to your body, and you can tell that the person affected you by sending a material substance your way, rather than just harming you directly with an arm motion, which would be known as “action at a distance.” It is not easy to do a similar observation to see whether there is some “stuff” that travels from the candle to your eye, or whether it is a case of action at a distance.

Newtonian physics includes both action at a distance (e.g., the earth’s gravitational force on a falling object) and contact forces such as the normal force, which only allow distant objects to exert forces on each other by shooting some substance across the space between them (e.g., a garden hose spraying out water that exerts a force on a bush).

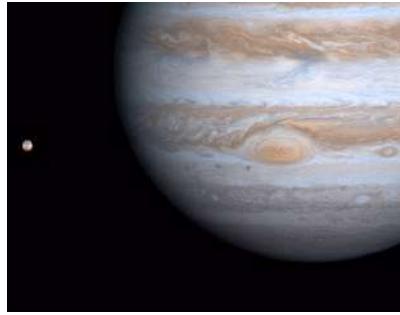
One piece of evidence that the candle sends out stuff that travels to your eye is that as in figure a, intervening transparent substances can make the candle appear to be in the wrong location, suggesting that light is a thing that can be bumped off course. Many people would dismiss this kind of observation as an optical illusion, however. (Some optical illusions are purely neurological or psychological effects, although some others, including this one, turn out to be caused by the behavior of light itself.)

A more convincing way to decide in which category light belongs is to find out if it takes time to get from the candle to your eye; in Newtonian physics, action at a distance is supposed to be instantaneous. The fact that we speak casually today of “the speed of light” implies that at some point in history, somebody succeeded in showing that light did not travel infinitely fast. Galileo tried, and failed, to detect a finite speed for light, by arranging with a person in a distant tower to signal back and forth with lanterns. Galileo uncovered his lantern, and when the other person saw the light, he uncovered his lantern. Galileo was unable to measure any time lag that was significant compared to the limitations of human reflexes.

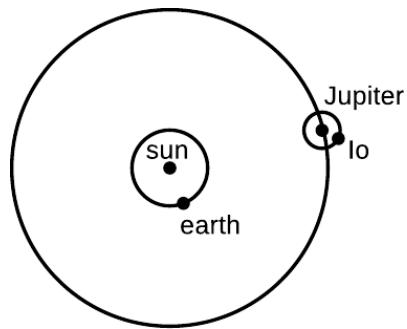


a / Light from a candle is bumped off course by a piece of glass. Inserting the glass causes the apparent location of the candle to shift. The same effect can be produced by taking off your eyeglasses and looking at which you see near the edge of the lens, but a flat piece of glass works just as well as a lens for this purpose.

The first person to prove that light's speed was finite, and to determine it numerically, was Ole Roemer, in a series of measurements around the year 1675. Roemer observed Io, one of Jupiter's moons, over a period of several years. Since Io presumably took the same amount of time to complete each orbit of Jupiter, it could be thought of as a very distant, very accurate clock. A practical and accurate pendulum clock had recently been invented, so Roemer could check whether the ratio of the two clocks' cycles, about 42.5 hours to 1 orbit, stayed exactly constant or changed a little. If the process of seeing the distant moon was instantaneous, there would be no reason for the two to get out of step. Even if the speed of light was finite, you might expect that the result would be only to offset one cycle relative to the other. The earth does not, however, stay at a constant distance from Jupiter and its moons. Since the distance is changing gradually due to the two planets' orbital motions, a finite speed of light would make the "Io clock" appear to run faster as the planets drew near each other, and more slowly as their separation increased. Roemer did find a variation in the apparent speed of Io's orbits, which caused Io's eclipses by Jupiter (the moments when Io passed in front of or behind Jupiter) to occur about 7 minutes early when the earth was closest to Jupiter, and 7 minutes late when it was farthest. Based on these measurements, Roemer estimated the speed of light to be approximately 2×10^8 m/s, which is in the right ballpark compared to modern measurements of 3×10^8 m/s. (I'm not sure whether the fairly large experimental error was mainly due to imprecise knowledge of the radius of the earth's orbit or limitations in the reliability of pendulum clocks.)



b / An image of Jupiter and its moon Io (left) from the Cassini probe.



c / The earth is moving toward Jupiter and Io. Since the distance is shrinking, it is taking less and less time for the light to get to us from Io, and Io appears to circle Jupiter more quickly than normal. Six months later, the earth will be on the opposite side of the sun, and receding from Jupiter and Io, so Io will appear to revolve around Jupiter more slowly.

Light can travel through a vacuum.

Many people are confused by the relationship between sound and light. Although we use different organs to sense them, there are some similarities. For instance, both light and sound are typically emitted in all directions by their sources. Musicians even use visual metaphors like "tone color," or "a bright timbre" to describe sound. One way to see that they are clearly different phenomena is to note their very different velocities. Sure, both are pretty fast compared to a flying arrow or a galloping horse, but as we have seen, the speed of light is so great as to appear instantaneous in most situations. The speed of sound, however, can easily be observed just by watching a group of schoolchildren a hundred feet away as they clap their hands to a song. There is an obvious delay between when you see their palms come together and when you hear the clap.

The fundamental distinction between sound and light is that sound is an oscillation in air pressure, so it requires air (or some other medium such as water) in which to travel. Today, we know that outer space is a vacuum, so the fact that we get light from the sun, moon and stars clearly shows that air is not necessary for the

propagation of light.

Discussion Questions

- A** If you observe thunder and lightning, you can tell how far away the storm is. Do you need to know the speed of sound, of light, or of both?
- B** When phenomena like X-rays and cosmic rays were first discovered, suggest a way one could have tested whether they were forms of light.
- C** Why did Roemer only need to know the radius of the earth's orbit, not Jupiter's, in order to find the speed of light?

12.1.2 Interaction of light with matter

Absorption of light

The reason why the sun feels warm on your skin is that the sunlight is being absorbed, and the light energy is being transformed into heat energy. The same happens with artificial light, so the net result of leaving a light turned on is to heat the room. It doesn't matter whether the source of the light is hot, like the sun, a flame, or an incandescent light bulb, or cool, like a fluorescent bulb. (If your house has electric heat, then there is absolutely no point in fastidiously turning off lights in the winter; the lights will help to heat the house at the same dollar rate as the electric heater.)

This process of heating by absorption is entirely different from heating by thermal conduction, as when an electric stove heats spaghetti sauce through a pan. Heat can only be conducted through matter, but there is vacuum between us and the sun, or between us and the filament of an incandescent bulb. Also, heat conduction can only transfer heat energy from a hotter object to a colder one, but a cool fluorescent bulb is perfectly capable of heating something that had already started out being warmer than the bulb itself.

How we see nonluminous objects

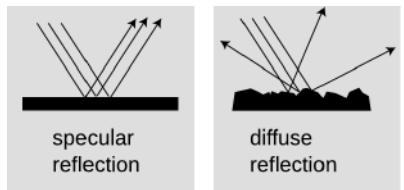
Not all the light energy that hits an object is transformed into heat. Some is reflected, and this leads us to the question of how we see nonluminous objects. If you ask the average person how we see a light bulb, the most likely answer is "The light bulb makes light, which hits our eyes." But if you ask how we see a book, they are likely to say "The bulb lights up the room, and that lets me see the book." All mention of light actually entering our eyes has mysteriously disappeared.

Most people would disagree if you told them that light was reflected from the book to the eye, because they think of reflection as something that mirrors do, not something that a book does. They associate reflection with the formation of a reflected image, which does not seem to appear in a piece of paper.

Imagine that you are looking at your reflection in a nice smooth piece of aluminum foil, fresh off the roll. You perceive a face, not a



d / Two self-portraits of the author, one taken in a mirror and one with a piece of aluminum foil.



e / Specular and diffuse reflection.

piece of metal. Perhaps you also see the bright reflection of a lamp over your shoulder behind you. Now imagine that the foil is just a little bit less smooth. The different parts of the image are now a little bit out of alignment with each other. Your brain can still recognize a face and a lamp, but it's a little scrambled, like a Picasso painting. Now suppose you use a piece of aluminum foil that has been crumpled up and then flattened out again. The parts of the image are so scrambled that you cannot recognize an image. Instead, your brain tells you you're looking at a rough, silvery surface.

Mirror-like reflection at a specific angle is known as specular reflection, and random reflection in many directions is called diffuse reflection. Diffuse reflection is how we see nonluminous objects. Specular reflection only allows us to see images of objects other than the one doing the reflecting. In top part of figure d, imagine that the rays of light are coming from the sun. If you are looking down at the reflecting surface, there is no way for your eye-brain system to tell that the rays are not really coming from a sun down below you.

Figure f shows another example of how we can't avoid the conclusion that light bounces off of things other than mirrors. The lamp is one I have in my house. It has a bright bulb, housed in a completely opaque bowl-shaped metal shade. The only way light can get out of the lamp is by going up out of the top of the bowl. The fact that I can read a book in the position shown in the figure means that light must be bouncing off of the ceiling, then bouncing off of the book, then finally getting to my eye.

This is where the shortcomings of the Greek theory of vision become glaringly obvious. In the Greek theory, the light from the bulb and my mysterious "eye rays" are both supposed to go to the book, where they collide, allowing me to see the book. But we now have a total of four objects: lamp, eye, book, and ceiling. Where does the ceiling come in? Does it also send out its own mysterious "ceiling rays," contributing to a three-way collision at the book? That would just be too bizarre to believe!

The differences among white, black, and the various shades of gray in between is a matter of what percentage of the light they absorb and what percentage they reflect. That's why light-colored clothing is more comfortable in the summer, and light-colored upholstery in a car stays cooler than dark upholstery.

Numerical measurement of the brightness of light

We have already seen that the physiological sensation of loudness relates to the sound's intensity (power per unit area), but is not directly proportional to it. If sound A has an intensity of 1 nW/m^2 , sound B is 10 nW/m^2 , and sound C is 100 nW/m^2 , then the increase in loudness from B to C is perceived to be the same as the increase from A to B, not ten times greater. That is, the sensation of loudness is logarithmic.

The same is true for the brightness of light. Brightness is related to power per unit area, but the psychological relationship is a logarithmic one rather than a proportionality. For doing physics, it's the power per unit area that we're interested in. The relevant unit is W/m^2 . One way to determine the brightness of light is to measure the increase in temperature of a black object exposed to the light. The light energy is being converted to heat energy, and the amount of heat energy absorbed in a given amount of time can be related to the power absorbed, using the known heat capacity of the object. More practical devices for measuring light intensity, such as the light meters built into some cameras, are based on the conversion of light into electrical energy, but these meters have to be calibrated somehow against heat measurements.

Discussion Questions

A The curtains in a room are drawn, but a small gap lets light through, illuminating a spot on the floor. It may or may not also be possible to see the beam of sunshine crossing the room, depending on the conditions. What's going on?

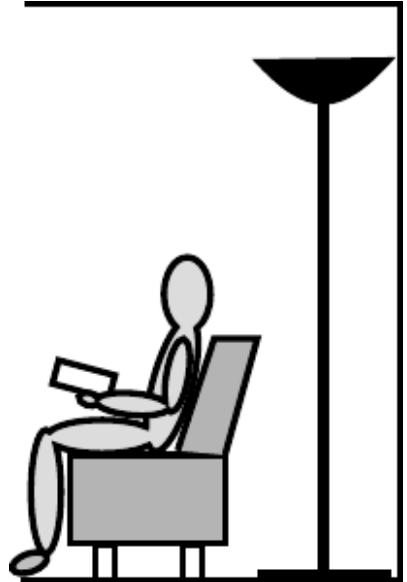
B Laser beams are made of light. In science fiction movies, laser beams are often shown as bright lines shooting out of a laser gun on a spaceship. Why is this scientifically incorrect?

C A documentary film-maker went to Harvard's 1987 graduation ceremony and asked the graduates, on camera, to explain the cause of the seasons. Only two out of 23 were able to give a correct explanation, but you now have all the information needed to figure it out for yourself, assuming you didn't already know. The figure shows the earth in its winter and summer positions relative to the sun. Hint: Consider the units used to measure the brightness of light, and recall that the sun is lower in the sky in winter, so its rays are coming in at a shallower angle.

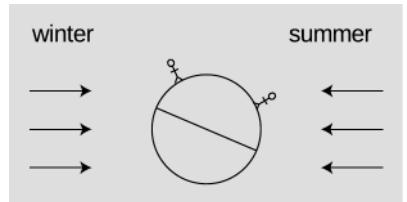
12.1.3 The ray model of light

Models of light

Note how I've been casually diagramming the motion of light with pictures showing light rays as lines on the page. More formally, this is known as the ray model of light. The ray model of light seems natural once we convince ourselves that light travels through space, and observe phenomena like sunbeams coming through holes in clouds. Having already been introduced to the concept of light as an electromagnetic wave, you know that the ray model is not the



f / Light bounces off of the ceiling, then off of the book.



g / Discussion question C.

ultimate truth about light, but the ray model is simpler, and in any case science always deals with models of reality, not the ultimate nature of reality. The following table summarizes three models of light.

h / Three models of light.

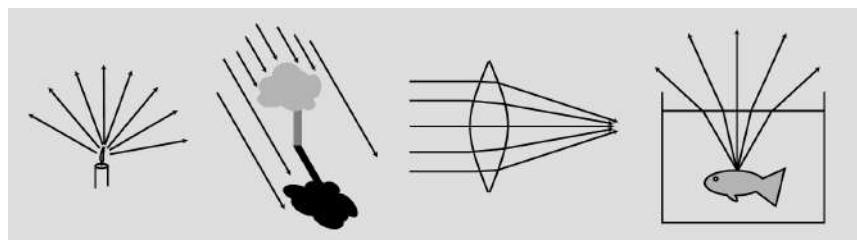
ray model		Advantage: Simplicity.
wave model		Advantage: Color is described naturally in terms of wavelength. Required in order to explain the interaction of light with material objects of sizes comparable to or smaller than a wavelength of light.
particle model		Required in order to explain the interaction of light with individual atoms. At the atomic level, it becomes apparent that a beam of light has a certain graininess to it.

The ray model is a generic one. By using it we can discuss the path taken by the light, without committing ourselves to any specific description of what it is that is moving along that path. We will use the nice simple ray model for most of our treatment of optics, and with it we can analyze a great many devices and phenomena. Not until section 12.5 will we concern ourselves specifically with wave optics, although in the intervening chapters I will sometimes analyze the same phenomenon using both the ray model and the wave model.

Note that the statements about the applicability of the various models are only rough guides. For instance, wave interference effects are often detectable, if small, when light passes around an obstacle that is quite a bit bigger than a wavelength. Also, the criterion for when we need the particle model really has more to do with energy scales than distance scales, although the two turn out to be related.

The alert reader may have noticed that the wave model is required at scales smaller than a wavelength of light (on the order of a micrometer for visible light), and the particle model is demanded on the atomic scale or lower (a typical atom being a nanometer or so in size). This implies that at the smallest scales we need *both* the wave model and the particle model. They appear incompatible, so how can we simultaneously use both? The answer is that they are not as incompatible as they seem. Light is both a wave and a particle,

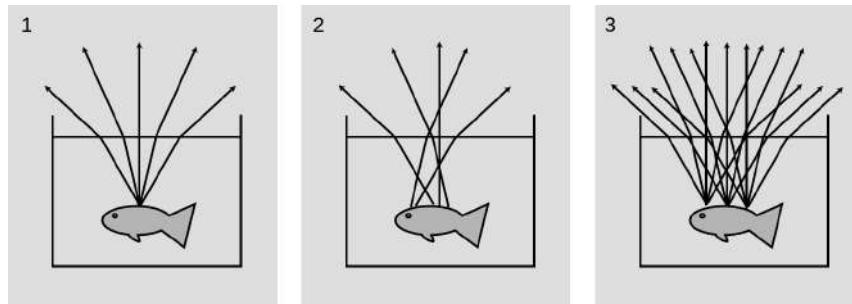
but a full understanding of this apparently nonsensical statement is a topic for section 13.2.



i / Examples of ray diagrams.

Ray diagrams

Without even knowing how to use the ray model to calculate anything numerically, we can learn a great deal by drawing ray diagrams. For instance, if you want to understand how eyeglasses help you to see in focus, a ray diagram is the right place to start. Many students under-utilize ray diagrams in optics and instead rely on rote memorization or plugging into formulas. The trouble with memorization and plug-ins is that they can obscure what's really going on, and it is easy to get them wrong. Often the best plan is to do a ray diagram first, then do a numerical calculation, then check that your numerical results are in reasonable agreement with what you expected from the ray diagram.



j / 1. Correct. 2. Incorrect: implies that diffuse reflection only gives one ray from each reflecting point. 3. Correct, but unnecessarily complicated

Figure j shows some guidelines for using ray diagrams effectively. The light rays bend when they pass out through the surface of the water (a phenomenon that we'll discuss in more detail later). The rays appear to have come from a point above the goldfish's actual location, an effect that is familiar to people who have tried spear-fishing.

- A stream of light is not really confined to a finite number of narrow lines. We just draw it that way. In j/1, it has been necessary to choose a finite number of rays to draw (five), rather than the theoretically infinite number of rays that will diverge from that point.

- There is a tendency to conceptualize rays incorrectly as objects. In his Optics, Newton goes out of his way to caution the reader against this, saying that some people “consider ... the refraction of ... rays to be the bending or breaking of them in their passing out of one medium into another.” But a ray is a record of the path traveled by light, not a physical thing that can be bent or broken.
- In theory, rays may continue infinitely far into the past and future, but we need to draw lines of finite length. In j/1, a judicious choice has been made as to where to begin and end the rays. There is no point in continuing the rays any farther than shown, because nothing new and exciting is going to happen to them. There is also no good reason to start them earlier, before being reflected by the fish, because the direction of the diffusely reflected rays is random anyway, and unrelated to the direction of the original, incoming ray.
- When representing diffuse reflection in a ray diagram, many students have a mental block against drawing many rays fanning out from the same point. Often, as in example j/2, the problem is the misconception that light can only be reflected in one direction from one point.
- Another difficulty associated with diffuse reflection, example j/3, is the tendency to think that in addition to drawing many rays coming out of one point, we should also be drawing many rays coming from many points. In j/1, drawing many rays coming out of one point gives useful information, telling us, for instance, that the fish can be seen from any angle. Drawing many sets of rays, as in j/3, does not give us any more useful information, and just clutters up the picture in this example. The only reason to draw sets of rays fanning out from more than one point would be if different things were happening to the different sets.

Discussion Question

A Suppose an intelligent tool-using fish is spear-hunting for humans. Draw a ray diagram to show how the fish has to correct its aim. Note that although the rays are now passing from the air to the water, the same rules apply: the rays are closer to being perpendicular to the surface when they are in the water, and rays that hit the air-water interface at a shallow angle are bent the most.

12.1.4 Geometry of specular reflection

To change the motion of a material object, we use a force. Is there any way to exert a force on a beam of light? Experiments show that electric and magnetic fields do not deflect light beams, so apparently light has no electric charge. Light also has no mass, so

until the twentieth century it was believed to be immune to gravity as well. Einstein predicted that light beams would be very slightly deflected by strong gravitational fields, and he was proved correct by observations of rays of starlight that came close to the sun, but obviously that's not what makes mirrors and lenses work!

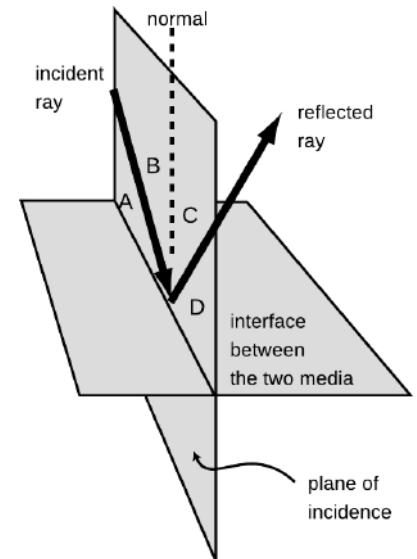
If we investigate how light is reflected by a mirror, we will find that the process is horrifically complex, but the final result is surprisingly simple. What actually happens is that the light is made of electric and magnetic fields, and these fields accelerate the electrons in the mirror. Energy from the light beam is momentarily transformed into extra kinetic energy of the electrons, but because the electrons are accelerating they re-radiate more light, converting their kinetic energy back into light energy. We might expect this to result in a very chaotic situation, but amazingly enough, the electrons move together to produce a new, reflected beam of light, which obeys two simple rules:

- The angle of the reflected ray is the same as that of the incident ray.
- The reflected ray lies in the plane containing the incident ray and the normal (perpendicular) line. This plane is known as the plane of incidence.

The two angles can be defined either with respect to the normal, like angles B and C in the figure, or with respect to the reflecting surface, like angles A and D. There is a convention of several hundred years' standing that one measures the angles with respect to the normal, but the rule about equal angles can logically be stated either as $B=C$ or as $A=D$.

The phenomenon of reflection occurs only at the boundary between two media, just like the change in the speed of light that passes from one medium to another. As we have seen in section 6.2, this is the way all waves behave.

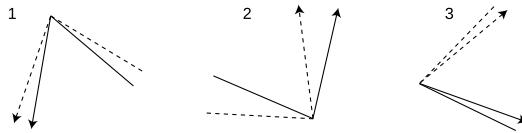
Most people are surprised by the fact that light can be reflected back from a less dense medium. For instance, if you are diving and you look up at the surface of the water, you will see a reflection of yourself.



k / The geometry of specular reflection.

self-check A

Each of these diagrams is supposed to show two different rays being reflected from the same point on the same mirror. Which are correct, and which are incorrect?



▷ Answer, p. 1065

Reversibility of light rays

The fact that specular reflection displays equal angles of incidence and reflection means that there is a symmetry: if the ray had come in from the right instead of the left in the figure above, the angles would have looked exactly the same. This is not just a pointless detail about specular reflection. It's a manifestation of a very deep and important fact about nature, which is that the laws of physics do not distinguish between past and future. Cannonballs and planets have trajectories that are equally natural in reverse, and so do light rays. This type of symmetry is called time-reversal symmetry.

Typically, time-reversal symmetry is a characteristic of any process that does not involve heat. For instance, the planets do not experience any friction as they travel through empty space, so there is no frictional heating. We should thus expect the time-reversed versions of their orbits to obey the laws of physics, which they do. In contrast, a book sliding across a table does generate heat from friction as it slows down, and it is therefore not surprising that this type of motion does not appear to obey time-reversal symmetry. A book lying still on a flat table is never observed to spontaneously start sliding, sucking up heat energy and transforming it into kinetic energy.

Similarly, the only situation we've observed so far where light does not obey time-reversal symmetry is absorption, which involves heat. Your skin absorbs visible light from the sun and heats up, but we never observe people's skin to glow, converting heat energy into visible light. People's skin does glow in infrared light, but that doesn't mean the situation is symmetric. Even if you absorb infrared, you don't emit visible light, because your skin isn't hot enough to glow in the visible spectrum.

These apparent heat-related asymmetries are not actual asymmetries in the laws of physics. The interested reader may wish to learn more about this from optional chapter 5 on thermodynamics.

Ray tracing on a computer

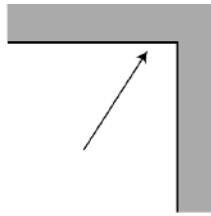
example 1

A number of techniques can be used for creating artificial visual scenes in computer graphics. Figure I shows such a scene, which was created by the brute-force technique of simply constructing a very detailed ray diagram on a computer. This technique requires a great deal of computation, and is therefore too slow to

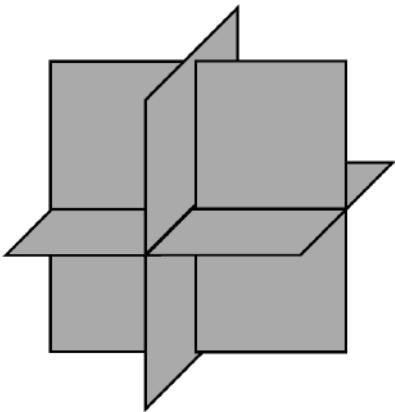
be used for video games and computer-animated movies. One trick for speeding up the computation is to exploit the reversibility of light rays. If one was to trace every ray emitted by every illuminated surface, only a tiny fraction of those would actually end up passing into the virtual “camera,” and therefore almost all of the computational effort would be wasted. One can instead start a ray at the camera, trace it backward in time, and see where it would have come from. With this technique, there is no wasted effort.



I / This photorealistic image of a nonexistent countertop was produced completely on a computer, by computing a complicated ray diagram.



m / Discussion question B.



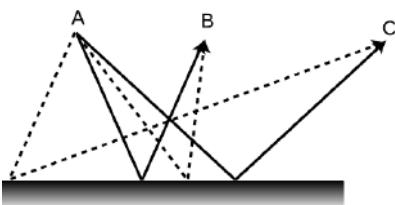
n / Discussion question C.

Discussion Questions

A If a light ray has a velocity vector with components c_x and c_y , what will happen when it is reflected from a surface that lies along the y axis? Make sure your answer does not imply a change in the ray's speed.

B Generalizing your reasoning from discussion question A, what will happen to the velocity components of a light ray that hits a corner, as shown in the figure, and undergoes two reflections?

C Three pieces of sheet metal arranged perpendicularly as shown in the figure form what is known as a radar corner. Let's assume that the radar corner is large compared to the wavelength of the radar waves, so that the ray model makes sense. If the radar corner is bathed in radar rays, at least some of them will undergo three reflections. Making a further generalization of your reasoning from the two preceding discussion questions, what will happen to the three velocity components of such a ray? What would the radar corner be useful for?



o / The solid lines are physically possible paths for light rays traveling from A to B and from A to C. They obey the principle of least time. The dashed lines do not obey the principle of least time, and are not physically possible.

12.1.5 * The principle of least time for reflection

We had to choose between an unwieldy explanation of reflection at the atomic level and a simpler geometric description that was not as fundamental. There is a third approach to describing the interaction of light and matter which is very deep and beautiful. Emphasized by the twentieth-century physicist Richard Feynman, it is called the principle of least time, or Fermat's principle.

Let's start with the motion of light that is not interacting with matter at all. In a vacuum, a light ray moves in a straight line. This can be rephrased as follows: of all the conceivable paths light could follow from P to Q, the only one that is physically possible is the path that takes the least time.

What about reflection? If light is going to go from one point to another, being reflected on the way, the quickest path is indeed the one with equal angles of incidence and reflection. If the starting and ending points are equally far from the reflecting surface, o, it's not hard to convince yourself that this is true, just based on symmetry. There is also a tricky and simple proof, shown in figure p, for the more general case where the points are at different distances from the surface.

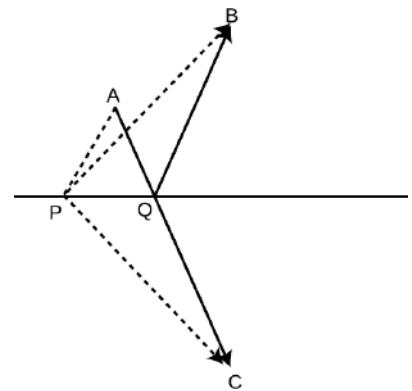
Not only does the principle of least time work for light in a

vacuum and light undergoing reflection, we will also see in a later chapter that it works for the bending of light when it passes from one medium into another.

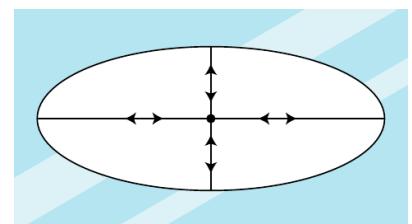
Although it is beautiful that the entire ray model of light can be reduced to one simple rule, the principle of least time, it may seem a little spooky to speak as if the ray of light is intelligent, and has carefully planned ahead to find the shortest route to its destination. How does it know in advance where it's going? What if we moved the mirror while the light was en route, so conditions along its planned path were not what it "expected?" The answer is that the principle of least time is really a shortcut for finding certain results of the wave model of light, which is the topic of the last chapter of this book.

There are a couple of subtle points about the principle of least time. First, the path does not have to be the quickest of all possible paths; it only needs to be quicker than any path that differs infinitesimally from it. In figure p, for instance, light could get from A to B either by the reflected path AQB or simply by going straight from A to B. Although AQB is not the shortest possible path, it cannot be shortened by changing it infinitesimally, e.g., by moving Q a little to the right or left. On the other hand, path APB is physically impossible, because it is possible to improve on it by moving point P infinitesimally to the right.

It's not quite right to call this the principle of *least* time. In figure q, for example, the four physically possible paths by which a ray can return to the center consist of two shortest-time paths and two longest-time paths. Strictly speaking, we should refer to the *principle of least or greatest time*, but most physicists omit the niceties, and assume that other physicists understand that both maxima and minima are possible.



p / Paths AQB and APB are two conceivable paths that a ray could follow to get from A to B with one reflection, but only AQB is physically possible. We wish to prove that the path AQB, with equal angles of incidence and reflection, is shorter than any other path, such as APB. The trick is to construct a third point, C, lying as far below the surface as B lies above it. Then path AQC is a straight line whose length is the same as AQB's, and path APC has the same length as path APB. Since AQC is straight, it must be shorter than any other path such as APC that connects A and C, and therefore AQB must be shorter than any path such as APB.



q / Light is emitted at the center of an elliptical mirror. There are four physically possible paths by which a ray can be reflected and return to the center.

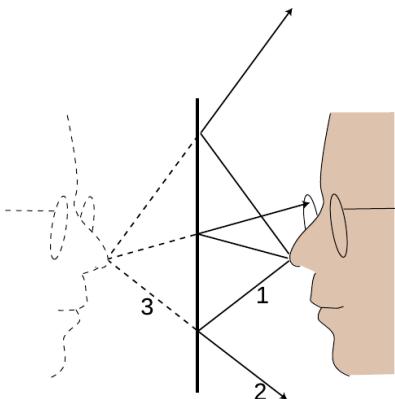
12.2 Images by reflection

Infants are always fascinated by the antics of the Baby in the Mirror. Now if you want to know something about mirror images that most people don't understand, try this. First bring this page closer and closer to your eyes, until you can no longer focus on it without straining. Then go in the bathroom and see how close you can get your face to the surface of the mirror before you can no longer easily focus on the image of your own eyes. You will find that the shortest comfortable eye-mirror distance is much less than the shortest comfortable eye-paper distance. This demonstrates that the image of your face in the mirror acts as if it had depth and existed in the space *behind* the mirror. If the image was like a flat picture in a book, then you wouldn't be able to focus on it from such a short distance.

In this chapter we will study the images formed by flat and curved mirrors on a qualitative, conceptual basis. Although this type of image is not as commonly encountered in everyday life as images formed by lenses, images formed by reflection are simpler to understand, so we discuss them first. In section 12.3 we will turn to a more mathematical treatment of images made by reflection. Surprisingly, the same equations can also be applied to lenses, which are the topic of section 12.4.

12.2.1 A virtual image

We can understand a mirror image using a ray diagram. Figure a shows several light rays, 1, that originated by diffuse reflection at the person's nose. They bounce off the mirror, producing new rays, 2. To anyone whose eye is in the right position to get one of these rays, they appear to have come from a behind the mirror, 3, where they would have originated from a single point. This point is where the tip of the image-person's nose appears to be. A similar analysis applies to every other point on the person's face, so it looks as though there was an entire face behind the mirror. The customary way of describing the situation requires some explanation:



a / An image formed by a mirror.

Customary description in physics: There is an image of the face behind the mirror.

Translation: The pattern of rays coming from the mirror is exactly the same as it would be if there were a face behind the mirror. Nothing is really behind the mirror.

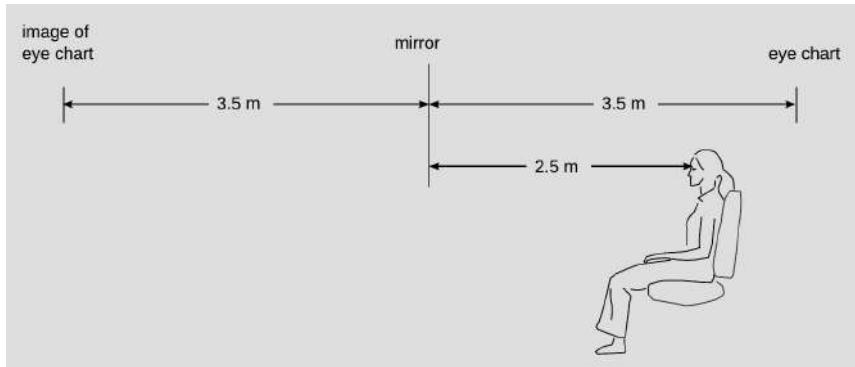
This is referred to as a *virtual* image, because the rays do not actually cross at the point behind the mirror. They only appear to have originated there.

self-check B

Imagine that the person in figure a moves his face down quite a bit — a couple of feet in real life, or a few inches on this scale drawing. The mirror stays where it is. Draw a new ray diagram. Will there still be an image? If so, where is it visible from?

▷ Answer, p. 1065

The geometry of specular reflection tells us that rays 1 and 2 are at equal angles to the normal (the imaginary perpendicular line piercing the mirror at the point of reflection). This means that ray 2's imaginary continuation, 3, forms the same angle with the mirror as ray 1. Since each ray of type 3 forms the same angles with the mirror as its partner of type 1, we see that the distance of the image from the mirror is the same as that of the actual face from the mirror, and it lies directly across from it. The image therefore appears to be the same size as the actual face.



b / Example 2.

An eye exam

example 2

Figure b shows a typical setup in an optometrist's examination room. The patient's vision is supposed to be tested at a distance of 6 meters (20 feet in the U.S.), but this distance is larger than the amount of space available in the room. Therefore a mirror is used to create an image of the eye chart behind the wall.

The Praxinoscope

example 3

Figure c shows an old-fashioned device called a praxinoscope, which displays an animated picture when spun. The removable strip of paper with the pictures printed on it has twice the radius of the inner circle made of flat mirrors, so each picture's virtual image is at the center. As the wheel spins, each picture's image is replaced by the next.



c / The praxinoscope.

Discussion Question

A The figure shows an object that is off to one side of a mirror. Draw a ray diagram. Is an image formed? If so, where is it, and from which directions would it be visible?

O



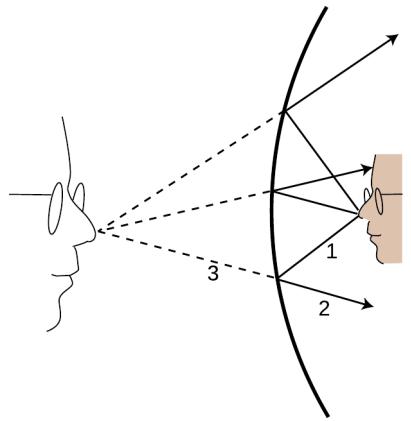
12.2.2 Curved mirrors

An image in a flat mirror is a pretechnological example: even animals can look at their reflections in a calm pond. We now pass to our first nontrivial example of the manipulation of an image by technology: an image in a curved mirror. Before we dive in, let's consider why this is an important example. If it was just a question of memorizing a bunch of facts about curved mirrors, then you would rightly rebel against an effort to spoil the beauty of your liberally educated brain by force-feeding you technological trivia. The reason this is an important example is not that curved mirrors are so important in and of themselves, but that the results we derive for curved bowl-shaped mirrors turn out to be true for a large class of other optical devices, including mirrors that bulge outward rather than inward, and lenses as well. A microscope or a telescope is simply a combination of lenses or mirrors or both. What you're really learning about here is the basic building block of all optical devices from movie projectors to octopus eyes.

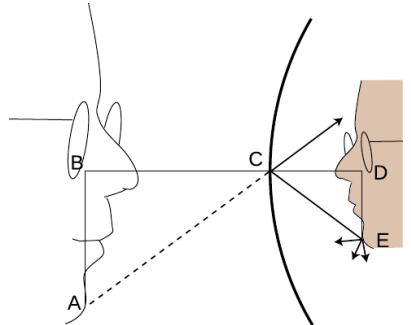
Because the mirror in figure d is curved, it bends the rays back closer together than a flat mirror would: we describe it as *converging*. Note that the term refers to what it does to the light rays, not to the physical shape of the mirror's surface. (The surface itself would be described as *concave*. The term is not all that hard to remember, because the hollowed-out interior of the mirror is like a cave.) It is surprising but true that all the rays like 3 really do converge on a point, forming a good image. We will not prove this fact, but it is true for any mirror whose curvature is gentle enough and that is symmetric with respect to rotation about the perpendicular line passing through its center (not asymmetric like a potato chip). The old-fashioned method of making mirrors and lenses is by grinding them in grit by hand, and this automatically tends to produce an almost perfect spherical surface.

Bending a ray like 2 inward implies bending its imaginary continuation 3 outward, in the same way that raising one end of a seesaw causes the other end to go down. The image therefore forms deeper behind the mirror. This doesn't just show that there is extra distance between the image-nose and the mirror; it also implies that the image itself is bigger from front to back. It has been *magnified* in the front-to-back direction.

It is easy to prove that the same magnification also applies to the image's other dimensions. Consider a point like E in figure e. The trick is that out of all the rays diffusely reflected by E, we pick the one that happens to head for the mirror's center, C. The equal-angle property of specular reflection plus a little straightforward geometry easily leads us to the conclusion that triangles ABC and CDE are the same shape, with ABC being simply a scaled-up version of CDE. The magnification of depth equals the ratio BC/CD, and the up-



d / An image formed by a curved mirror.



e / The image is magnified by the same factor in depth and in its other dimensions.



f / Increased magnification always comes at the expense of decreased field of view.

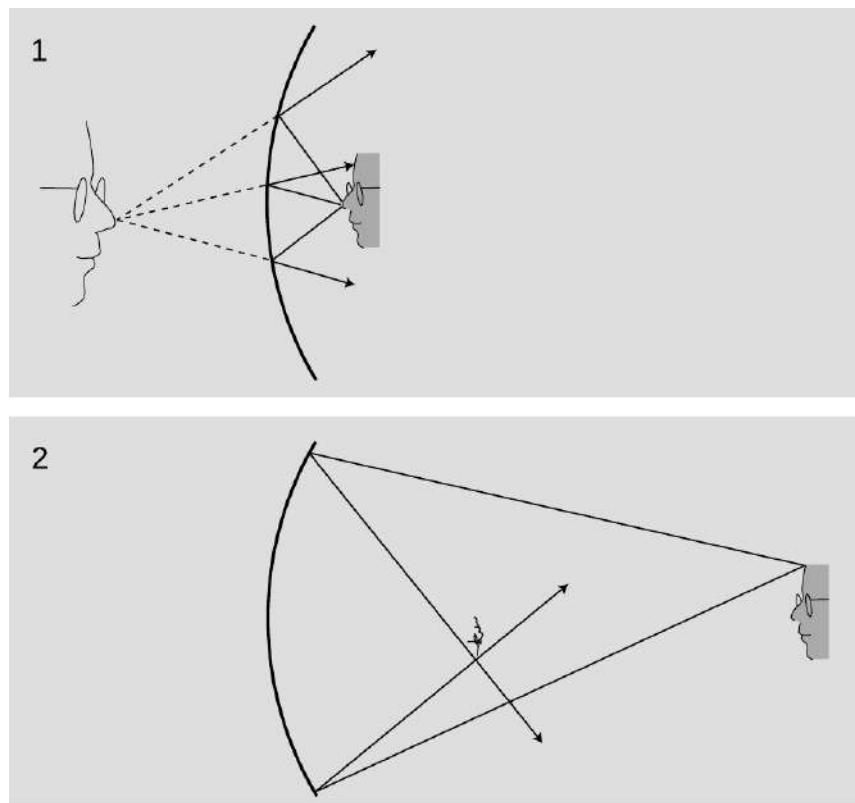
down magnification is AB/DE . A repetition of the same proof shows that the magnification in the third dimension (out of the page) is also the same. This means that the image-head is simply a larger version of the real one, without any distortion. The scaling factor is called the magnification, M . The image in the figure is magnified by a factor $M = 1.9$.

Note that we did not explicitly specify whether the mirror was a sphere, a paraboloid, or some other shape. However, we assumed that a focused image would be formed, which would not necessarily be true, for instance, for a mirror that was asymmetric or very deeply curved.

12.2.3 A real image

If we start by placing an object very close to the mirror, $g/1$, and then move it farther and farther away, the image at first behaves as we would expect from our everyday experience with flat mirrors, receding deeper and deeper behind the mirror. At a certain point, however, a dramatic change occurs. When the object is more than a certain distance from the mirror, $g/2$, the image appears upside-down and in *front* of the mirror.

g / 1. A virtual image. **2.** A real image. As you'll verify in homework problem 12, the image is upside-down



Here's what's happened. The mirror bends light rays inward, but when the object is very close to it, as in $g/1$, the rays coming from a

given point on the object are too strongly diverging (spreading) for the mirror to bring them back together. On reflection, the rays are still diverging, just not as strongly diverging. But when the object is sufficiently far away, $g/2$, the mirror is only intercepting the rays that came out in a narrow cone, and it is able to bend these enough so that they will reconverge.

Note that the rays shown in the figure, which both originated at the same point on the object, reunite when they cross. The point where they cross is the image of the point on the original object. This type of image is called a *real image*, in contradistinction to the virtual images we've studied before.

Definition: A real image is one where rays actually cross. A virtual image is a point from which rays only appear to have come.

The use of the word "real" is perhaps unfortunate. It sounds as though we are saying the image was an actual material object, which of course it is not.

The distinction between a real image and a virtual image is an important one, because a real image can be projected onto a screen or photographic film. If a piece of paper is inserted in figure $g/2$ at the location of the image, the image will be visible on the paper (provided the object is bright and the room is dark). Your eye uses a lens to make a real image on the retina.

self-check C

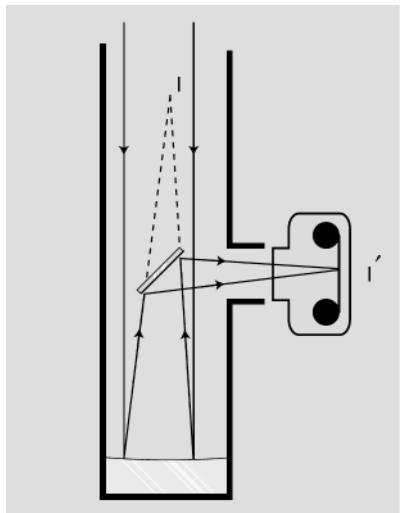
Sketch another copy of the face in figure $g/1$, even farther from the mirror, and draw a ray diagram. What has happened to the location of the image?

▷ Answer, p. 1066

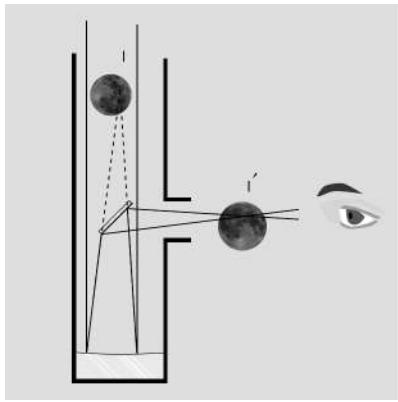
12.2.4 Images of images

If you are wearing glasses right now, then the light rays from the page are being manipulated first by your glasses and then by the lens of your eye. You might think that it would be extremely difficult to analyze this, but in fact it is quite easy. In any series of optical elements (mirrors or lenses or both), each element works on the rays furnished by the previous element in exactly the same manner as if the image formed by the previous element was an actual object.

Figure h shows an example involving only mirrors. The Newtonian telescope, invented by Isaac Newton, consists of a large curved mirror, plus a second, flat mirror that brings the light out of the tube. (In very large telescopes, there may be enough room to put a camera or even a person inside the tube, in which case the second mirror is not needed.) The tube of the telescope is not vital; it is mainly a structural element, although it can also be helpful for blocking out stray light. The lens has been removed from the front of the camera body, and is not needed for this setup. Note that the



h / A Newtonian telescope being used with a camera.

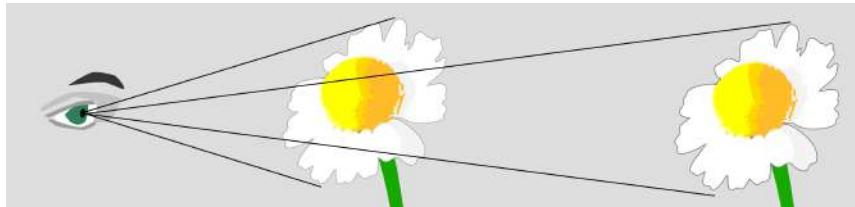


i / A Newtonian telescope being used for visual rather than photographic observing. In real life, an eyepiece lens is normally used for additional magnification, but this simpler setup will also work.

two sample rays have been drawn parallel, because an astronomical telescope is used for viewing objects that are extremely far away. These two “parallel” lines actually meet at a certain point, say a crater on the moon, so they can’t actually be perfectly parallel, but they are parallel for all practical purposes since we would have to follow them upward for a quarter of a million miles to get to the point where they intersect.

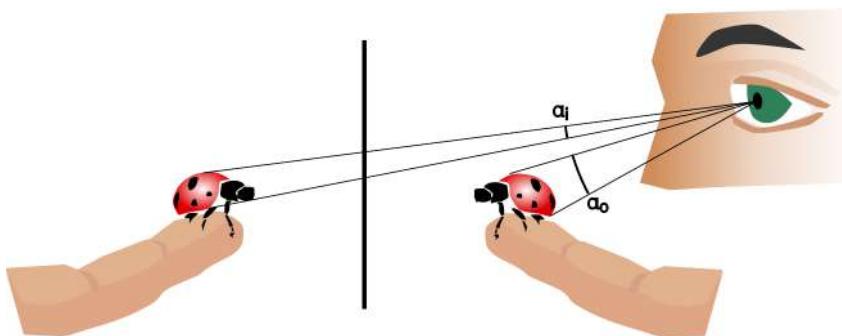
The large curved mirror by itself would form an image I , but the small flat mirror creates an image of the image, I' . The relationship between I and I' is exactly the same as it would be if I was an actual object rather than an image: I and I' are at equal distances from the plane of the mirror, and the line between them is perpendicular to the plane of the mirror.

One surprising wrinkle is that whereas a flat mirror used by itself forms a virtual image of an object that is real, here the mirror is forming a real image of virtual image I . This shows how pointless it would be to try to memorize lists of facts about what kinds of images are formed by various optical elements under various circumstances. You are better off simply drawing a ray diagram.



j / The angular size of the flower depends on its distance from the eye.

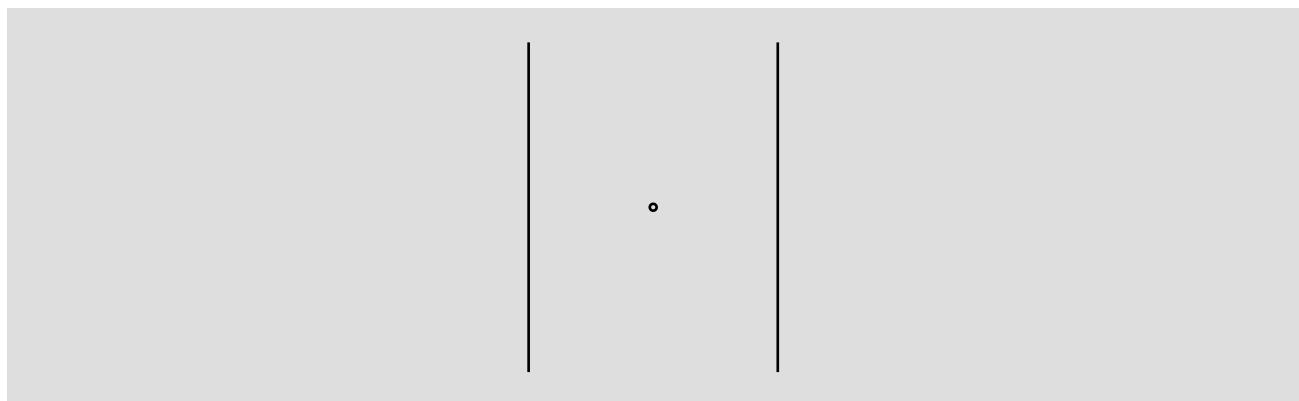
Although the main point here was to give an example of an image of an image, figure i also shows an interesting case where we need to make the distinction between *magnification* and *angular magnification*. If you are looking at the moon through this telescope, then the images I and I' are much *smaller* than the actual moon. Otherwise, for example, image I would not fit inside the telescope! However, these images are very close to your eye compared to the actual moon. The small size of the image has been more than compensated for by the shorter distance. The important thing here is the amount of *angle* within your field of view that the image covers, and it is this angle that has been increased. The factor by which it is increased is called the *angular magnification*, M_a .



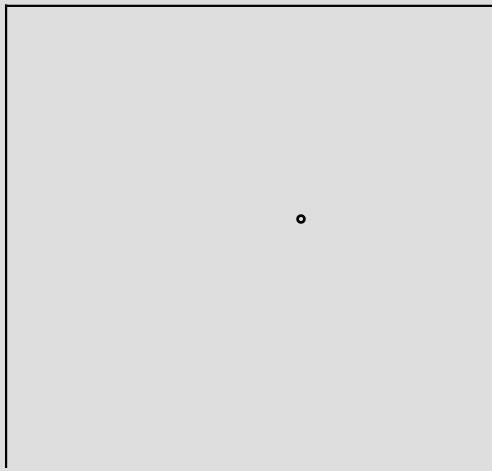
k / The person uses a mirror to get a view of both sides of the ladybug. Although the flat mirror has $M = 1$, it doesn't give an angular magnification of 1. The image is farther from the eye than the object, so the angular magnification $M_a = \alpha_i/\alpha_o$ is less than one.

Discussion Questions

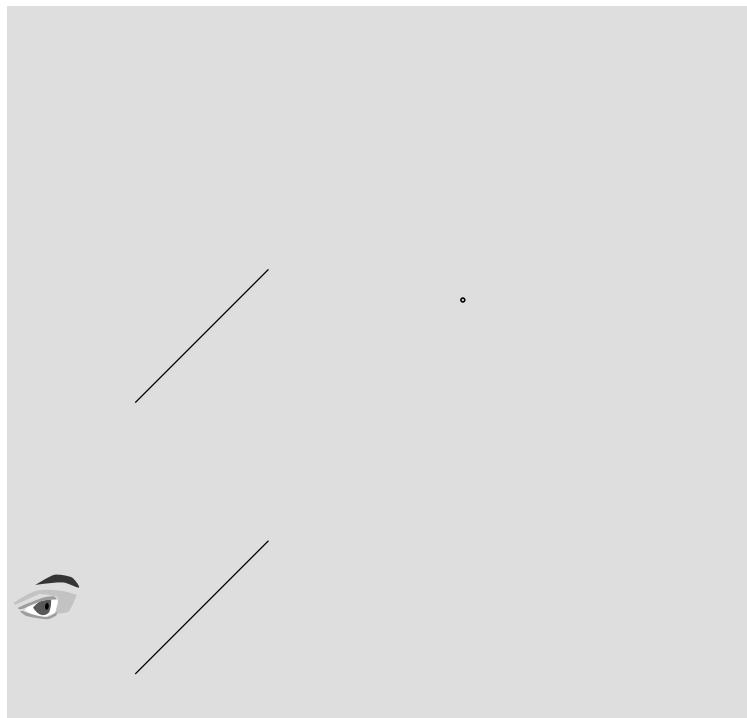
- A** Locate the images of you that will be formed if you stand between two parallel mirrors.



B Locate the images formed by two perpendicular mirrors, as in the figure. What happens if the mirrors are not perfectly perpendicular?

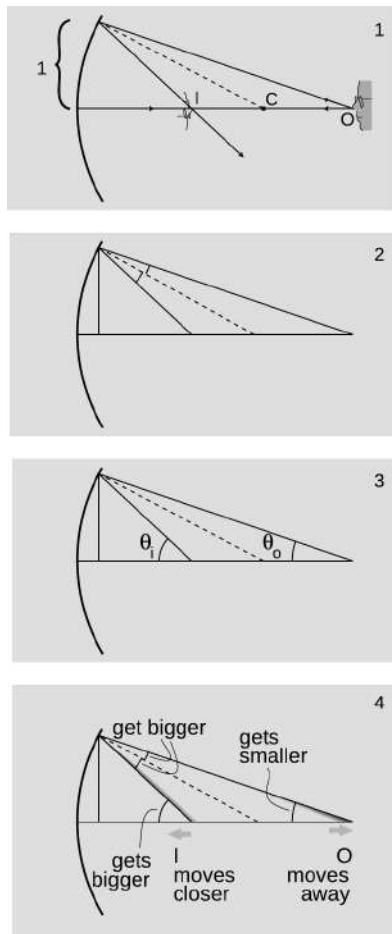


C Locate the images formed by the periscope.



12.3 Images, quantitatively

It sounds a bit odd when a scientist refers to a theory as “beautiful,” but to those in the know it makes perfect sense. One mark of a beautiful theory is that it surprises us by being simple. The mathematical theory of lenses and curved mirrors gives us just such a surprise. We expect the subject to be complex because there are so many cases: a converging mirror forming a real image, a diverging lens that makes a virtual image, and so on for a total of six possibilities. If we want to predict the location of the images in all these situations, we might expect to need six different equations, and six more for predicting magnifications. Instead, it turns out that we can use just one equation for the location of the image and one equation for its magnification, and these two equations work in all the different cases with no changes except for plus and minus signs. This is the kind of thing the physicist Eugene Wigner referred to as “the unreasonable effectiveness of mathematics.” Sometimes we can find a deeper reason for this kind of unexpected simplicity, but sometimes it almost seems as if God went out of Her way to make the secrets of universe susceptible to attack by the human thought-tool called math.



a / The relationship between the object's position and the image's can be expressed in terms of the angles θ_o and θ_i .

12.3.1 A real image formed by a converging mirror

Location of the image

We will now derive the equation for the location of a real image formed by a converging mirror. We assume for simplicity that the mirror is spherical, but actually this isn't a restrictive assumption, because any shallow, symmetric curve can be approximated by a sphere. The shape of the mirror can be specified by giving the location of its center, C. A deeply curved mirror is a sphere with a small radius, so C is close to it, while a weakly curved mirror has C farther away. Given the point O where the object is, we wish to find the point I where the image will be formed.

To locate an image, we need to track a minimum of two rays coming from the same point. Since we have proved in the previous chapter that this type of image is not distorted, we can use an on-axis point, O, on the object, as in figure a/1. The results we derive will also hold for off-axis points, since otherwise the image would have to be distorted, which we know is not true. We let one of the rays be the one that is emitted along the axis; this ray is especially easy to trace, because it bounces straight back along the axis again. As our second ray, we choose one that strikes the mirror at a distance of 1 from the axis. “One what?” asks the astute reader. The answer is that it doesn't really matter. When a mirror has shallow curvature, all the reflected rays hit the same point, so 1 could be expressed in any units you like. It could, for instance, be 1 cm, unless your mirror is smaller than 1 cm!

The only way to find out anything mathematical about the rays is to use the sole mathematical fact we possess concerning specular reflection: the incident and reflected rays form equal angles with respect to the normal, which is shown as a dashed line. Therefore the two angles shown in figure a/2 are the same, and skipping some straightforward geometry, this leads to the visually reasonable result that the two angles in figure a/3 are related as follows:

$$\theta_i + \theta_o = \text{constant}$$

(Note that θ_i and θ_o , which are measured from the image and the object, not from the eye like the angles we referred to in discussing angular magnification on page 786.) For example, move O farther from the mirror. The top angle in figure a/2 is increased, so the bottom angle must increase by the same amount, causing the image point, I, to move closer to the mirror. In terms of the angles shown in figure a/3, the more distant object has resulted in a smaller angle θ_o , while the closer image corresponds to a larger θ_i ; One angle increases by the same amount that the other decreases, so their sum remains constant. These changes are summarized in figure a/4.

The sum $\theta_i + \theta_o$ is a constant. What does this constant represent? Geometrically, we interpret it as double the angle made by the dashed radius line. Optically, it is a measure of the strength of the mirror, i.e., how strongly the mirror focuses light, and so we call it the focal angle, θ_f ,

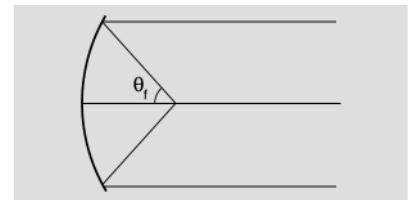
$$\theta_i + \theta_o = \theta_f.$$

Suppose, for example, that we wish to use a quick and dirty optical test to determine how strong a particular mirror is. We can lay it on the floor as shown in figure c, and use it to make an image of a lamp mounted on the ceiling overhead, which we assume is very far away compared to the radius of curvature of the mirror, so that the mirror intercepts only a very narrow cone of rays from the lamp. This cone is so narrow that its rays are nearly parallel, and θ_o is nearly zero. The real image can be observed on a piece of paper. By moving the paper nearer and farther, we can bring the image into focus, at which point we know the paper is located at the image point. Since $\theta_o \approx 0$, we have $\theta_i \approx \theta_f$, and we can then determine this mirror's focal angle either by measuring θ_i directly with a protractor, or indirectly via trigonometry. A strong mirror will bring the rays together to form an image close to the mirror, and these rays will form a blunt-angled cone with a large θ_i and θ_f .

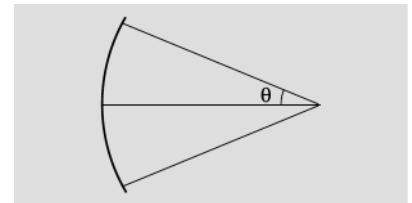
An alternative optical test

example 4

- ▷ Figure c shows an alternative optical test. Rather than placing the object at infinity as in figure b, we adjust it so that the image is right on top of the object. Points O and I coincide, and the rays are reflected right back on top of themselves. If we measure the angle θ shown in figure c, how can we find the focal angle?



b / The geometrical interpretation of the focal angle.



c / Example 4, an alternative test for finding the focal angle. The mirror is the same as in figure b.

▷ The object and image angles are the same; the angle labeled θ in the figure equals both of them. We therefore have $\theta_i + \theta_o = \theta = \theta_f$. Comparing figures b and c, it is indeed plausible that the angles are related by a factor of two.

At this point, we could consider our work to be done. Typically, we know the strength of the mirror, and we want to find the image location for a given object location. Given the mirror's focal angle and the object location, we can determine θ_o by trigonometry, subtract to find $\theta_i = \theta_f - \theta_o$, and then do more trig to find the image location.

There is, however, a shortcut that can save us from doing so much work. Figure a/3 shows two right triangles whose legs of length 1 coincide and whose acute angles are θ_o and θ_i . These can be related by trigonometry to the object and image distances shown in figure d:

$$\tan \theta_o = 1/d_o \quad \tan \theta_i = 1/d_i$$

Ever since section 12.2, we've been assuming small angles. For small angles, we can use the small-angle approximation $\tan x \approx x$ (for x in radians), giving simply

$$\theta_o = 1/d_o \quad \theta_i = 1/d_i.$$

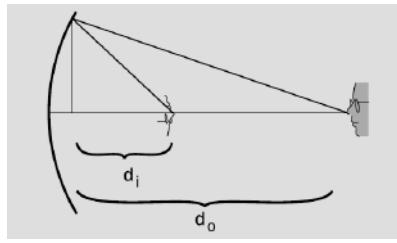
We likewise define a distance called the focal length, f according to $\theta_f = 1/f$. In figure b, f is the distance from the mirror to the place where the rays cross. We can now reexpress the equation relating the object and image positions as

$$\frac{1}{f} = \frac{1}{d_i} + \frac{1}{d_o}.$$

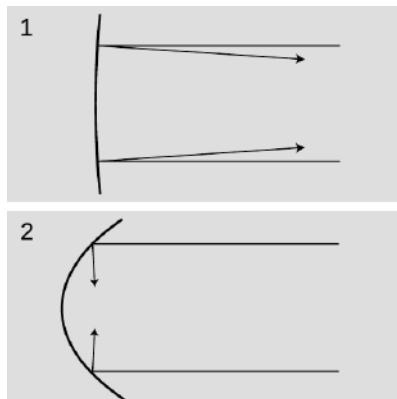
Figure e summarizes the interpretation of the focal length and focal angle.¹

Which form is better, $\theta_f = \theta_i + \theta_o$ or $1/f = 1/d_i + 1/d_o$? The angular form has in its favor its simplicity and its straightforward visual interpretation, but there are two reasons why we might prefer the second version. First, the numerical values of the angles depend on what we mean by "one unit" for the distance shown as 1 in

¹There is a standard piece of terminology which is that the "focal point" is the point lying on the optical axis at a distance from the mirror equal to the focal length. This term isn't particularly helpful, because it names a location where nothing normally happens. In particular, it is *not* normally the place where the rays come to a focus! — that would be the *image* point. In other words, we don't normally have $d_i = f$, unless perhaps $d_o = \infty$. A recent online discussion among some physics teachers (<https://carnot.physics.buffalo.edu/archives>, Feb. 2006) showed that many disliked the terminology, felt it was misleading, or didn't know it and would have misinterpreted it if they had come across it. That is, it appears to be what grammarians call a "skunked term" — a word that bothers half the population when it's used incorrectly, and the other half when it's used correctly.



d / The object and image distances



e / Mirror 1 is weaker than mirror 2. It has a shallower curvature, a longer focal length, and a smaller focal angle. It reflects rays at angles not much different than those that would be produced with a flat mirror.

figure a/1. Second, it is usually easier to measure distances rather than angles, so the distance form is more convenient for number crunching. Neither form is superior overall, and we will often need to use both to solve any given problem.²

A searchlight

example 5

Suppose we need to create a parallel beam of light, as in a searchlight. Where should we place the lightbulb? A parallel beam has zero angle between its rays, so $\theta_i = 0$. To place the lightbulb correctly, however, we need to know a distance, not an angle: the distance d_o between the bulb and the mirror. The problem involves a mixture of distances and angles, so we need to get everything in terms of one or the other in order to solve it. Since the goal is to find a distance, let's figure out the image distance corresponding to the given angle $\theta_i = 0$. These are related by $d_i = 1/\theta_i$, so we have $d_i = \infty$. (Yes, dividing by zero gives infinity. Don't be afraid of infinity. Infinity is a useful problem-solving device.) Solving the distance equation for d_o , we have

$$\begin{aligned} d_o &= (1/f - 1/d_i)^{-1} \\ &= (1/f - 0)^{-1} \\ &= f \end{aligned}$$

The bulb has to be placed at a distance from the mirror equal to its focal point.

Diopters

example 6

An equation like $d_i = 1/\theta_i$ really doesn't make sense in terms of units. Angles are unitless, since radians aren't really units, so the right-hand side is unitless. We can't have a left-hand side with units of distance if the right-hand side of the same equation is unitless. This is an artifact of my cavalier statement that the conical bundles of rays spread out to a distance of 1 from the axis where they strike the mirror, without specifying the units used to measure this 1. In real life, optometrists define the thing we're calling $\theta_i = 1/d_i$ as the "dioptic strength" of a lens or mirror, and measure it in units of inverse meters (m^{-1}), also known as diopters ($1 D=1 m^{-1}$).

Magnification

We have already discussed in the previous chapter how to find the magnification of a virtual image made by a curved mirror. The result is the same for a real image, and we omit the proof, which is very similar. In our new notation, the result is $M = d_i/d_o$. A numerical example is given in subsection 12.3.2.

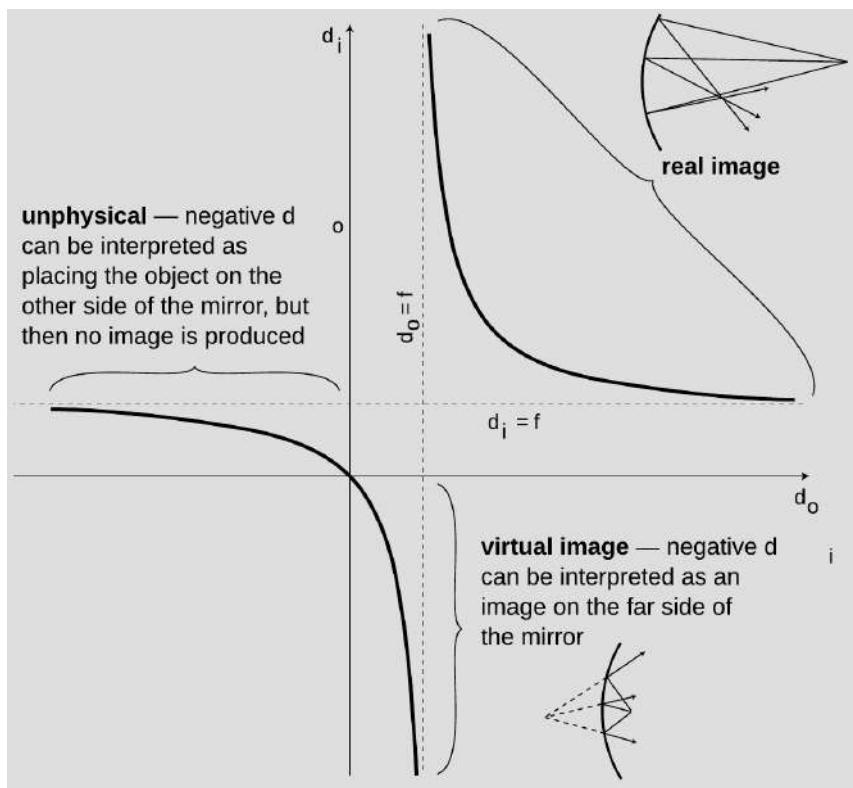
²I would like to thank Fouad Ajami for pointing out the pedagogical advantages of using both equations side by side.

12.3.2 Other cases with curved mirrors

The equation $d_i = (1/f - 1/d_o)^{-1}$ can easily produce a negative result, but we have been thinking of d_i as a distance, and distances can't be negative. A similar problem occurs with $\theta_i = \theta_f - \theta_o$ for $\theta_o > \theta_f$. What's going on here?

The interpretation of the angular equation is straightforward. As we bring the object closer and closer to the image, θ_o gets bigger and bigger, and eventually we reach a point where $\theta_o = \theta_f$ and $\theta_i = 0$. This large object angle represents a bundle of rays forming a cone that is very broad, so broad that the mirror can no longer bend them back so that they reconverge on the axis. The image angle $\theta_i = 0$ represents an outgoing bundle of rays that are parallel. The outgoing rays never cross, so this is not a real image, unless we want to be charitable and say that the rays cross at infinity. If we go on bringing the object even closer, we get a virtual image.

f / A graph of the image distance d_i as a function of the object distance d_o .



To analyze the distance equation, let's look at a graph of d_i as a function of d_o . The branch on the upper right corresponds to the case of a real image. Strictly speaking, this is the only part of the graph that we've proven corresponds to reality, since we never did any geometry for other cases, such as virtual images. As discussed in the previous section, making d_o bigger causes d_i to become smaller, and vice-versa.

Letting d_o be less than f is equivalent to $\theta_o > \theta_f$: a virtual image is produced on the far side of the mirror. This is the first example of Wigner’s “unreasonable effectiveness of mathematics” that we have encountered in optics. Even though our proof depended on the assumption that the image was real, the equation we derived turns out to be applicable to virtual images, provided that we either interpret the positive and negative signs in a certain way, or else modify the equation to have different positive and negative signs.

self-check D

Interpret the three places where, in physically realistic parts of the graph, the graph approaches one of the dashed lines. [This will come more naturally if you have learned the concept of limits in a math class.] ▷

Answer, p. 1066

A flat mirror

example 7

We can even apply the equation to a flat mirror. As a sphere gets bigger and bigger, its surface is more and more gently curved. The planet Earth is so large, for example, that we cannot even perceive the curvature of its surface. To represent a flat mirror, we let the mirror’s radius of curvature, and its focal length, become infinite. Dividing by infinity gives zero, so we have

$$1/d_o = -1/d_i,$$

or

$$d_o = -d_i.$$

If we interpret the minus sign as indicating a virtual image on the far side of the mirror from the object, this makes sense.

It turns out that for any of the six possible combinations of real or virtual images formed by converging or diverging lenses or mirrors, we can apply equations of the form

$$\theta_f = \theta_i + \theta_o$$

and

$$\frac{1}{f} = \frac{1}{d_i} + \frac{1}{d_o},$$

with only a modification of plus or minus signs. There are two possible approaches here. The approach we have been using so far is the more popular approach in American textbooks: leave the equation the same, but attach interpretations to the resulting negative or positive values of the variables. The trouble with this approach is that one is then forced to memorize tables of sign conventions, e.g., that the value of d_i should be negative when the image is a virtual image formed by a converging mirror. Positive and negative

signs also have to be memorized for focal lengths. Ugh! It's highly unlikely that any student has ever retained these lengthy tables in his or her mind for more than five minutes after handing in the final exam in a physics course. Of course one can always look such things up when they are needed, but the effect is to turn the whole thing into an exercise in blindly plugging numbers into formulas.

As you have gathered by now, there is another method which I think is better, and which I'll use throughout the rest of this book. In this method, all distances and angles are *positive by definition*, and we put in positive and negative signs in the *equations* depending on the situation. (I thought I was the first to invent this method, but I've been told that this is known as the European sign convention, and that it's fairly common in Europe.) Rather than memorizing these signs, we start with the generic equations

$$\theta_f = \pm\theta_i \pm \theta_o$$

$$\frac{1}{f} = \pm\frac{1}{d_i} \pm \frac{1}{d_o},$$

and then determine the signs by a two-step method that depends on ray diagrams. There are really only two signs to determine, not four; the signs in the two equations match up in the way you'd expect. The method is as follows:

1. Use ray diagrams to decide whether θ_o and θ_i vary in the same way or in opposite ways. (In other words, decide whether making θ_o greater results in a greater value of θ_i or a smaller one.) Based on this, decide whether the two signs in the angle equation are the same or opposite. If the signs are opposite, go on to step 2 to determine which is positive and which is negative.
2. If the signs are opposite, we need to decide which is the positive one and which is the negative. Since the focal angle is never negative, the smaller angle must be the one with a minus sign.

In step 1, many students have trouble drawing the ray diagram correctly. For simplicity, you should always do your diagram for a point on the object that is on the axis of the mirror, and let one of your rays be the one that is emitted along the axis and reflected straight back on itself, as in the figures in subsection 12.3.1. As shown in figure a/4 in subsection 12.3.1, there are four angles involved: two at the mirror, one at the object (θ_o), and one at the image (θ_i). Make sure to draw in the normal to the mirror so that you can see the two angles at the mirror. These two angles are equal, so as you change the object position, they fan out or fan in, like opening or closing a book. Once you've drawn this effect, you should easily be able to tell whether θ_o and θ_i change in the same way or in opposite ways.

Although focal lengths are always positive in the method used in this book, you should be aware that diverging mirrors and lenses

are assigned negative focal lengths in the other method, so if you see a lens labeled $f = -30$ cm, you'll know what it means.

An anti-shoplifting mirror

example 8

Convenience stores often install a diverging mirror so that the clerk has a view of the whole store and can catch shoplifters. Use a ray diagram to show that the image is reduced, bringing more into the clerk's field of view. If the focal length of the mirror is 3.0 m, and the mirror is 7.0 m from the farthest wall, how deep is the image of the store?

As shown in ray diagram g/1, d_i is less than d_o . The magnification, $M = d_i/d_o$, will be less than one, i.e., the image is actually reduced rather than magnified.

Apply the method outlined above for determining the plus and minus signs. Step 1: The object is the point on the opposite wall. As an experiment, g/2, move the object closer. I did these drawings using illustration software, but if you were doing them by hand, you'd want to make the scale much larger for greater accuracy. Also, although I split figure g into two separate drawings in order to make them easier to understand, you're less likely to make a mistake if you do them on top of each other.

The two angles at the mirror fan out from the normal. Increasing θ_o has clearly made θ_i larger as well. (All four angles got bigger.) There must be a cancellation of the effects of changing the two terms on the right in the same way, and the only way to get such a cancellation is if the two terms in the angle equation have opposite signs:

$$\theta_f = +\theta_i - \theta_o$$

or

$$\theta_f = -\theta_i + \theta_o.$$

Step 2: Now which is the positive term and which is negative? Since the image angle is bigger than the object angle, the angle equation must be

$$\theta_f = \theta_i - \theta_o,$$

in order to give a positive result for the focal angle. The signs of the distance equation behave the same way:

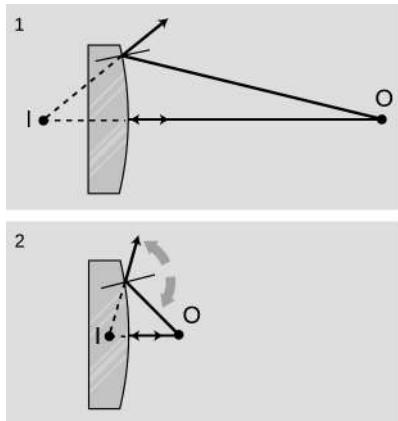
$$\frac{1}{f} = \frac{1}{d_i} - \frac{1}{d_o}.$$

Solving for d_i , we find

$$d_i = \left(\frac{1}{f} + \frac{1}{d_o} \right)^{-1}$$

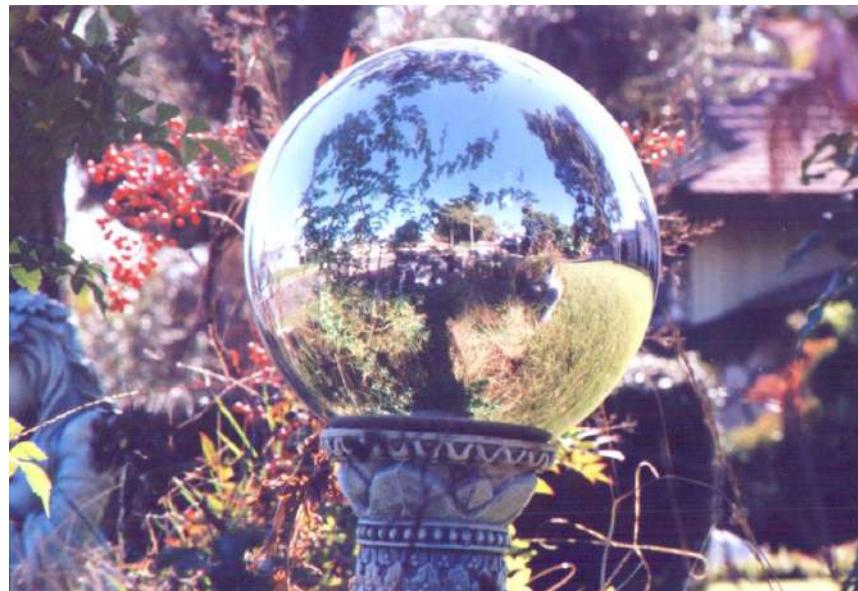
$$= 2.1 \text{ m.}$$

The image of the store is reduced by a factor of $2.1/7.0 = 0.3$, i.e., it is smaller by 70%.



g / Example 8.

h / A diverging mirror in the shape of a sphere. The image is reduced ($M < 1$). This is similar to example 8, but here the mirror's curve is not shallow.



A shortcut for real images

example 9

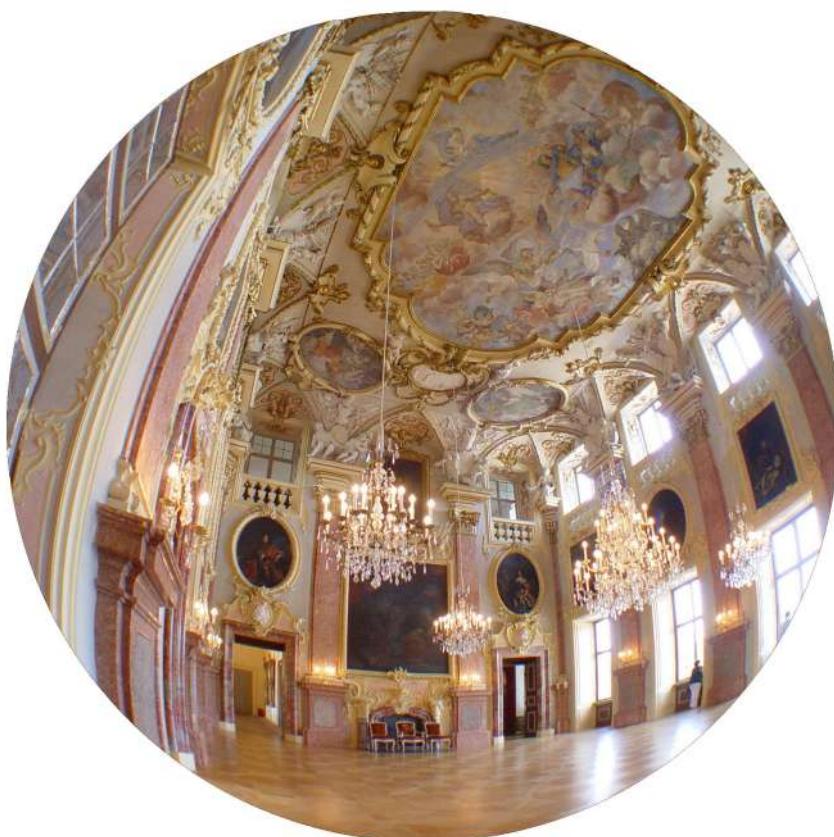
In the case of a real image, there is a shortcut for step 1, the determination of the signs. In a real image, the rays cross at both the object and the image. We can therefore time-reverse the ray diagram, so that all the rays are coming from the image and reconverging at the object. Object and image swap roles. Due to this time-reversal symmetry, the object and image cannot be treated differently in any of the equations, and they must therefore have the same signs. They are both positive, since they must add up to a positive result.

12.3.3 ★ Aberrations

An imperfection or distortion in an image is called an aberration. An aberration can be produced by a flaw in a lens or mirror, but even with a perfect optical surface some degree of aberration is unavoidable. To see why, consider the mathematical approximation we've been making, which is that the depth of the mirror's curve is small compared to d_o and d_i . Since only a flat mirror can satisfy this shallow-mirror condition perfectly, any curved mirror will deviate somewhat from the mathematical behavior we derived by assuming that condition. There are two main types of aberration in curved mirrors, and these also occur with lenses.

(1) An object on the axis of the lens or mirror may be imaged correctly, but off-axis objects may be out of focus or distorted. In a camera, this type of aberration would show up as a fuzziness or warping near the sides of the picture when the center was perfectly focused. An example of this is shown in figure i, and in that particular example, the aberration is not a sign that the equipment was of low quality or wasn't right for the job but rather an inevitable result of trying to flatten a panoramic view; in the limit of a 360-

degree panorama, the problem would be similar to the problem of representing the Earth's surface on a flat map, which can't be accomplished without distortion.



i / This photo was taken using a "fish-eye lens," which gives an extremely large field of view.

(2) The image may be sharp when the object is at certain distances and blurry when it is at other distances. The blurriness occurs because the rays do not all cross at exactly the same point. If we know in advance the distance of the objects with which the mirror or lens will be used, then we can optimize the shape of the optical surface to make in-focus images in that situation. For instance, a spherical mirror will produce a perfect image of an object that is at the center of the sphere, because each ray is reflected directly onto the radius along which it was emitted. For objects at greater distances, however, the focus will be somewhat blurry. In astronomy the objects being used are always at infinity, so a spherical mirror is a poor choice for a telescope. A different shape (a parabola) is better specialized for astronomy.

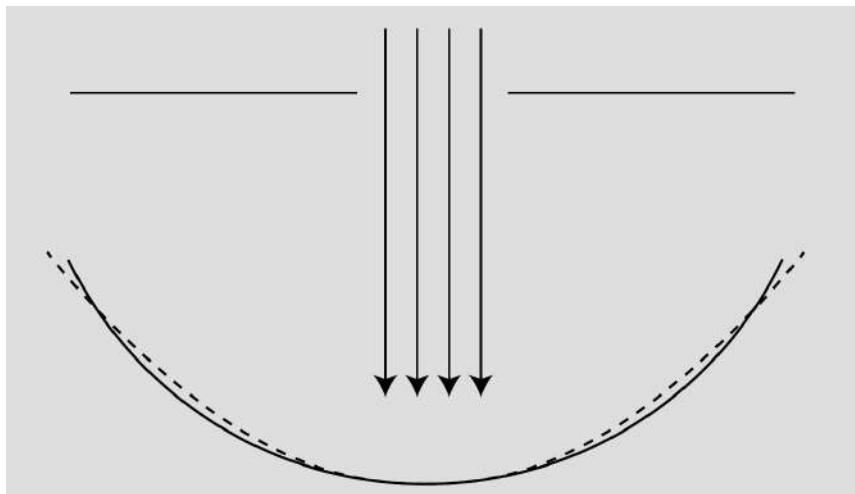
One way of decreasing aberration is to use a small-diameter mirror or lens, or block most of the light with an opaque screen with a hole in it, so that only light that comes in close to the axis can get

j / Spherical mirrors are the cheapest to make, but parabolic mirrors are better for making images of objects at infinity. A sphere has equal curvature everywhere, but a parabola has tighter curvature at its center and gentler curvature at the sides.



through. Either way, we are using a smaller portion of the lens or mirror whose curvature will be more shallow, thereby making the shallow-mirror (or thin-lens) approximation more accurate. Your eye does this by narrowing down the pupil to a smaller hole. In a camera, there is either an automatic or manual adjustment, and narrowing the opening is called “stopping down.” The disadvantage of stopping down is that light is wasted, so the image will be dimmer or a longer exposure must be used.

k / Even though the spherical mirror (solid line) is not well adapted for viewing an object at infinity, we can improve its performance greatly by stopping it down. Now the only part of the mirror being used is the central portion, where its shape is virtually indistinguishable from a parabola (dashed line).



What I would suggest you take away from this discussion for the sake of your general scientific education is simply an understanding of what an aberration is, why it occurs, and how it can be reduced, not detailed facts about specific types of aberrations.

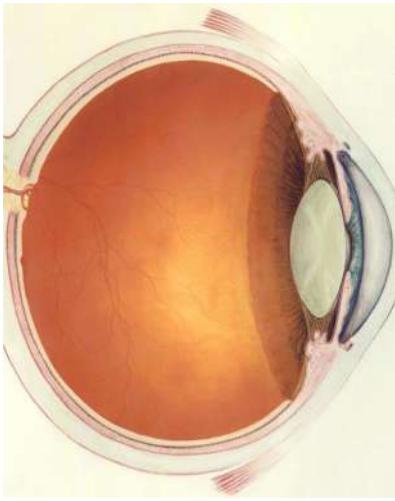


I / The Hubble Space Telescope was placed into orbit with faulty optics in 1990. Its main mirror was supposed to have been nearly parabolic, since it is an astronomical telescope, meant for producing images of objects at infinity. However, contractor Perkin Elmer had delivered a faulty mirror, which produced aberrations. The large photo shows astronauts putting correcting mirrors in place in 1993. The two small photos show images produced by the telescope before and after the fix.

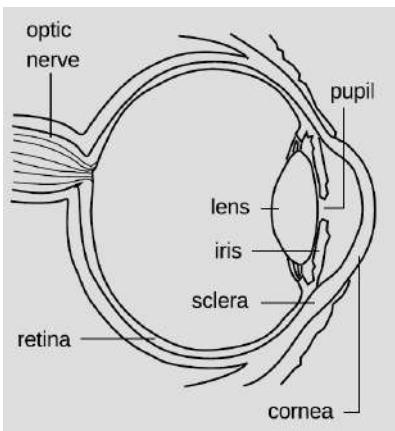
12.4 Refraction

Economists normally consider free markets to be the natural way of judging the monetary value of something, but social scientists also use questionnaires to gauge the relative value of privileges, disadvantages, or possessions that cannot be bought or sold. They ask people to *imagine* that they could trade one thing for another and ask which they would choose. One interesting result is that the average light-skinned person in the U.S. would rather lose an arm than suffer the racist treatment routinely endured by African-Americans. Even more impressive is the value of sight. Many prospective parents can imagine without too much fear having a deaf child, but would have a far more difficult time coping with raising a blind one.

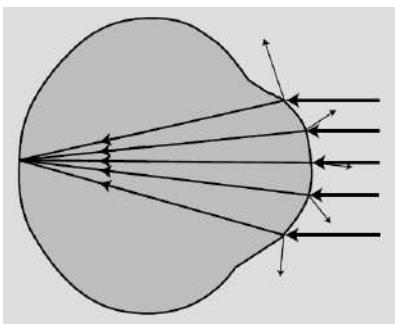
So great is the value attached to sight that some have imbued it with mystical aspects. Joan of Arc saw visions, and my college has a “vision statement.” Christian fundamentalists who perceive a conflict between evolution and their religion have claimed that the eye is such a perfect device that it could never have arisen through a process as helter-skelter as evolution, or that it could not have evolved because half of an eye would be useless. In fact, the structure of an eye is fundamentally dictated by physics, and it has arisen separately by evolution somewhere between eight and 40 times, depending on which biologist you ask. We humans have a version of the eye that can be traced back to the evolution of a light-sensitive “eye spot” on the head of an ancient invertebrate. A sunken pit then developed so that the eye would only receive light from one direction, allowing the organism to tell where the light was coming from. (Modern flatworms have this type of eye.) The top of the pit then became partially covered, leaving a hole, for even greater directionality (as in the nautilus). At some point the cavity became filled with jelly, and this jelly finally became a lens, resulting in the



a / A human eye.



b / The anatomy of the eye.



c / A simplified optical diagram of the eye. Light rays are bent when they cross from the air into the eye. (A little of the incident rays' energy goes into the reflected rays rather than the ones transmitted into the eye.)

general type of eye that we share with the bony fishes and other vertebrates. Far from being a perfect device, the vertebrate eye is marred by a serious design flaw due to the lack of planning or intelligent design in evolution: the nerve cells of the retina and the blood vessels that serve them are all in front of the light-sensitive cells, blocking part of the light. Squids and other molluscs, whose eyes evolved on a separate branch of the evolutionary tree, have a more sensible arrangement, with the light-sensitive cells out in front.

12.4.1 Refraction

Refraction

The fundamental physical phenomenon at work in the eye is that when light crosses a boundary between two media (such as air and the eye's jelly), part of its energy is reflected, but part passes into the new medium. In the ray model of light, we describe the original ray as splitting into a reflected ray and a transmitted one (the one that gets through the boundary). Of course the reflected ray goes in a direction that is different from that of the original one, according to the rules of reflection we have already studied. More surprisingly — and this is the crucial point for making your eye focus light — the transmitted ray is bent somewhat as well. This bending phenomenon is called *refraction*. The origin of the word is the same as that of the word “fracture,” i.e., the ray is bent or “broken.” (Keep in mind, however, that light rays are not physical objects that can really be “broken.”) Refraction occurs with all waves, not just light waves.

The actual anatomy of the eye, b, is quite complex, but in essence it is very much like every other optical device based on refraction. The rays are bent when they pass through the front surface of the eye, c. Rays that enter farther from the central axis are bent more, with the result that an image is formed on the retina. There is only one slightly novel aspect of the situation. In most human-built optical devices, such as a movie projector, the light is bent as it passes into a lens, bent again as it reemerges, and then reaches a focus beyond the lens. In the eye, however, the “screen” is inside the eye, so the rays are only refracted once, on entering the jelly, and never emerge again.

A common misconception is that the “lens” of the eye is what does the focusing. All the transparent parts of the eye are made of fairly similar stuff, so the dramatic change in medium is when a ray crosses from the air into the eye (at the outside surface of the cornea). This is where nearly all the refraction takes place. The lens medium differs only slightly in its optical properties from the rest of the eye, so very little refraction occurs as light enters and exits the lens. The lens, whose shape is adjusted by muscles attached to it, is only meant for fine-tuning the focus to form images of near or far objects.

Refractive properties of media

What are the rules governing refraction? The first thing to observe is that just as with reflection, the new, bent part of the ray lies in the same plane as the normal (perpendicular) and the incident ray, d.

If you try shooting a beam of light at the boundary between two substances, say water and air, you'll find that regardless of the angle at which you send in the beam, the part of the beam in the water is always closer to the normal line, e. It doesn't matter if the ray is entering the water or leaving, so refraction is symmetric with respect to time-reversal, f.

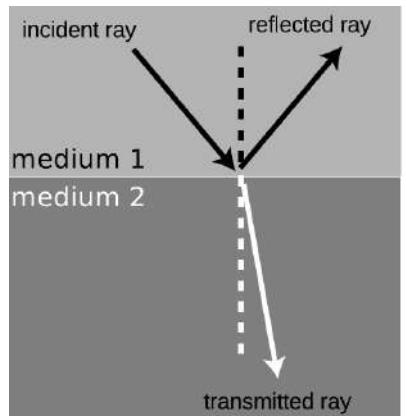
If, instead of water and air, you try another combination of substances, say plastic and gasoline, again you'll find that the ray's angle with respect to the normal is consistently smaller in one and larger in the other. Also, we find that if substance A has rays closer to normal than in B, and B has rays closer to normal than in C, then A has rays closer to normal than C. This means that we can rank-order all materials according to their refractive properties. Isaac Newton did so, including in his list many amusing substances, such as "Danzig vitriol" and "a pseudo-topazius, being a natural, pellucid, brittle, hairy stone, of a yellow color." Several general rules can be inferred from such a list:

- Vacuum lies at one end of the list. In refraction across the interface between vacuum and any other medium, the other medium has rays closer to the normal.
- Among gases, the ray gets closer to the normal if you increase the density of the gas by pressurizing it more.
- The refractive properties of liquid mixtures and solutions vary in a smooth and systematic manner as the proportions of the mixture are changed.
- Denser substances usually, but not always, have rays closer to the normal.

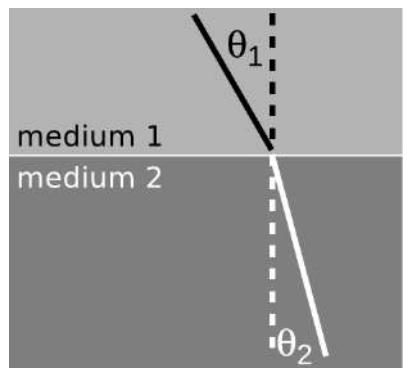
The second and third rules provide us with a method for measuring the density of an unknown sample of gas, or the concentration of a solution. The latter technique is very commonly used, and the CRC Handbook of Physics and Chemistry, for instance, contains extensive tables of the refractive properties of sugar solutions, cat urine, and so on.

Snell's law

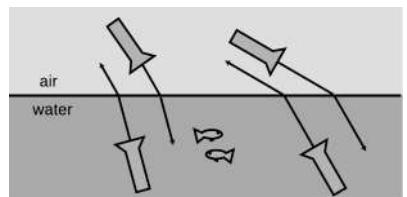
The numerical rule governing refraction was discovered by Snell, who must have collected experimental data something like what is



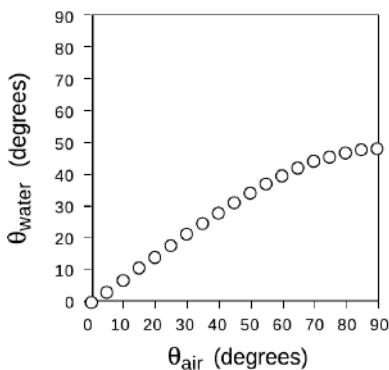
d / The incident, reflected, and transmitted (refracted) rays all lie in a plane that includes the normal (dashed line).



e / The angles θ_1 and θ_2 are related to each other, and also depend on the properties of the two media. Because refraction is time-reversal symmetric, there is no need to label the rays with arrowheads.



f / Refraction has time-reversal symmetry. Regardless of whether the light is going into or out of the water, the relationship between the two angles is the same, and the ray is closer to the normal while in the water.



g / The relationship between the angles in refraction.

shown on this graph and then attempted by trial and error to find the right equation. The equation he came up with was

$$\frac{\sin \theta_1}{\sin \theta_2} = \text{constant.}$$

The value of the constant would depend on the combination of media used. For instance, any one of the data points in the graph would have sufficed to show that the constant was 1.3 for an air-water interface (taking air to be substance 1 and water to be substance 2).

Snell further found that if media A and B gave a constant K_{AB} and media B and C gave a constant K_{BC} , then refraction at an interface between A and C would be described by a constant equal to the product, $K_{AC} = K_{AB}K_{BC}$. This is exactly what one would expect if the constant depended on the ratio of some number characterizing one medium to the number characteristic of the second medium. This number is called the *index of refraction* of the medium, written as n in equations. Since measuring the angles would only allow him to determine the *ratio* of the indices of refraction of two media, Snell had to pick some medium and define it as having $n = 1$. He chose to define vacuum as having $n = 1$. (The index of refraction of air at normal atmospheric pressure is 1.0003, so for most purposes it is a good approximation to assume that air has $n = 1$.) He also had to decide which way to define the ratio, and he chose to define it so that media with their rays closer to the normal would have larger indices of refraction. This had the advantage that denser media would typically have higher indices of refraction, and for this reason the index of refraction is also referred to as the optical density. Written in terms of indices of refraction, Snell's equation becomes

$$\frac{\sin \theta_1}{\sin \theta_2} = \frac{n_2}{n_1},$$

but rewriting it in the form

$$n_1 \sin \theta_1 = n_2 \sin \theta_2$$

[relationship between angles of rays at the interface between media with indices of refraction n_1 and n_2 ; angles are defined with respect to the normal]

makes us less likely to get the 1's and 2's mixed up, so this is the way most people remember Snell's law. A few indices of refraction are given in the back of the book.

self-check E

(1) What would the graph look like for two substances with the same index of refraction?

(2) Based on the graph, when does refraction at an air-water interface change the direction of a ray most strongly? ▷ Answer, p. 1066

Finding an angle using Snell's law

example 10

▷ A submarine shines its searchlight up toward the surface of the water. What is the angle α shown in the figure?

▷ The tricky part is that Snell's law refers to the angles with respect to the normal. Forgetting this is a very common mistake. The beam is at an angle of 30° with respect to the normal in the water. Let's refer to the air as medium 1 and the water as 2. Solving Snell's law for θ_1 , we find

$$\theta_1 = \sin^{-1} \left(\frac{n_2}{n_1} \sin \theta_2 \right).$$

As mentioned above, air has an index of refraction very close to 1, and water's is about 1.3, so we find $\theta_1 = 40^\circ$. The angle α is therefore 50° .

The index of refraction is related to the speed of light.

What neither Snell nor Newton knew was that there is a very simple interpretation of the index of refraction. This may come as a relief to the reader who is taken aback by the complex reasoning involving proportionalities that led to its definition. Later experiments showed that the index of refraction of a medium was inversely proportional to the speed of light in that medium. Since c is defined as the speed of light in vacuum, and $n = 1$ is defined as the index of refraction of vacuum, we have

$$n = \frac{c}{v}.$$

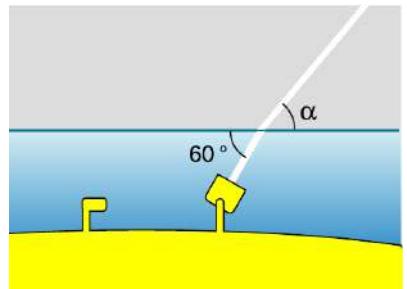
[n = medium's index of refraction, v = speed of light in that medium, c = speed of light in a vacuum]

Many textbooks start with this as the definition of the index of refraction, although that approach makes the quantity's name somewhat of a mystery, and leaves students wondering why c/v was used rather than v/c . It should also be noted that measuring angles of refraction is a far more practical method for determining n than direct measurement of the speed of light in the substance of interest.

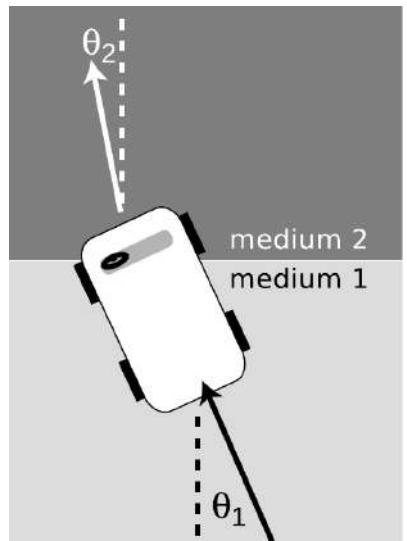
A mechanical model of Snell's law

Why should refraction be related to the speed of light? The mechanical model shown in the figure may help to make this more plausible. Suppose medium 2 is thick, sticky mud, which slows down the car. The car's right wheel hits the mud first, causing the right side of the car to slow down. This will cause the car to turn to the right until it moves far enough forward for the left wheel to cross into the mud. After that, the two sides of the car will once again be moving at the same speed, and the car will go straight.

Of course, light isn't a car. Why should a beam of light have anything resembling a "left wheel" and "right wheel?" After all,

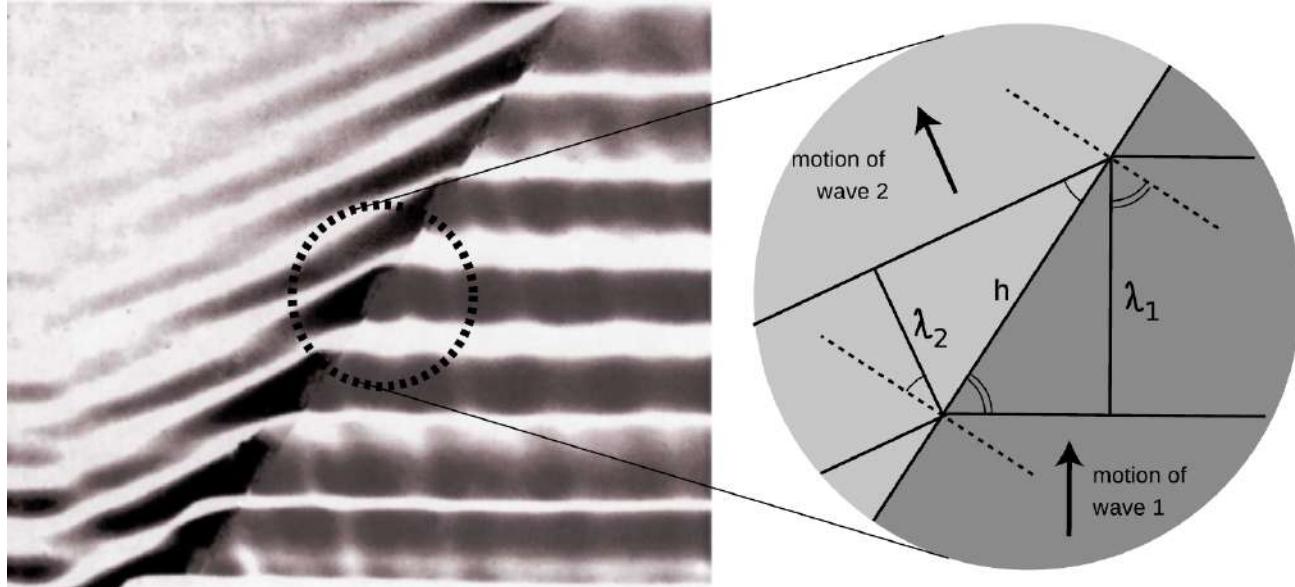


h / Example 10.



i / A mechanical model of refraction.

the mechanical model would predict that a motorcycle would go straight, and a motorcycle seems like a better approximation to a ray of light than a car. The whole thing is just a model, not a description of physical reality.



j / A derivation of Snell's law.

A derivation of Snell's law

However intuitively appealing the mechanical model may be, light is a wave, and we should be using wave models to describe refraction. In fact Snell's law can be derived quite simply from wave concepts. Figure j shows the refraction of a water wave. The water in the upper left part of the tank is shallower, so the speed of the waves is slower there, and their wavelengths is shorter. The reflected part of the wave is also very faintly visible.

In the close-up view on the right, the dashed lines are normals to the interface. The two marked angles on the right side are both equal to θ_1 , and the two on the left to θ_2 .

Trigonometry gives

$$\begin{aligned}\sin \theta_1 &= \lambda_1/h && \text{and} \\ \sin \theta_2 &= \lambda_2/h.\end{aligned}$$

Eliminating h by dividing the equations, we find

$$\frac{\sin \theta_1}{\sin \theta_2} = \frac{\lambda_1}{\lambda_2}.$$

The frequencies of the two waves must be equal or else they would get out of step, so by $v = f\lambda$ we know that their wavelengths are

proportional to their velocities. Combining $\lambda \propto v$ with $v \propto 1/n$ gives $\lambda \propto 1/n$, so we find

$$\frac{\sin \theta_1}{\sin \theta_2} = \frac{n_2}{n_1},$$

which is one form of Snell's law.

Ocean waves near and far from shore

example 11

Ocean waves are formed by winds, typically on the open sea, and the wavefronts are perpendicular to the direction of the wind that formed them. At the beach, however, you have undoubtedly observed that waves tend come in with their wavefronts very nearly (but not exactly) parallel to the shoreline. This is because the speed of water waves in shallow water depends on depth: the shallower the water, the slower the wave. Although the change from the fast-wave region to the slow-wave region is gradual rather than abrupt, there is still refraction, and the wave motion is nearly perpendicular to the normal in the slow region.

Color and refraction

In general, the speed of light in a medium depends both on the medium and on the wavelength of the light. Another way of saying it is that a medium's index of refraction varies with wavelength. This is why a prism can be used to split up a beam of white light into a rainbow. Each wavelength of light is refracted through a different angle.

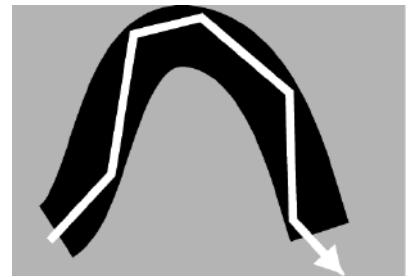
How much light is reflected, and how much is transmitted?

In section 6.2 we developed an equation for the percentage of the wave energy that is transmitted and the percentage reflected at a boundary between media. This was only done in the case of waves in one dimension, however, and rather than discuss the full three dimensional generalization it will be more useful to go into some qualitative observations about what happens. First, reflection happens only at the interface between two media, and two media with the same index of refraction act as if they were a single medium. Thus, at the interface between media with the same index of refraction, there is no reflection, and the ray keeps going straight. Continuing this line of thought, it is not surprising that we observe very little reflection at an interface between media with similar indices of refraction.

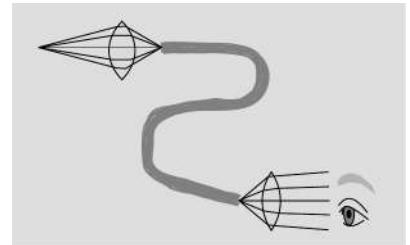
The next thing to note is that it is possible to have situations where no possible angle for the refracted ray can satisfy Snell's law. Solving Snell's law for θ_2 , we find

$$\theta_2 = \sin^{-1} \left(\frac{n_1}{n_2} \sin \theta_1 \right),$$

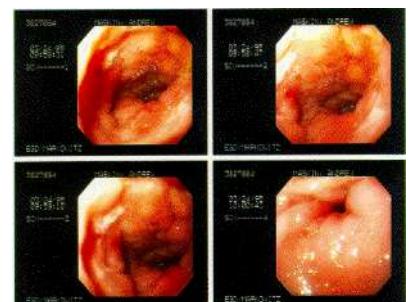
and if n_1 is greater than n_2 , then there will be large values of θ_1 for which the quantity $(n_1/n_2) \sin \theta$ is greater than one, meaning



k / Total internal reflection in a fiber-optic cable.



l / A simplified drawing of a surgical endoscope. The first lens forms a real image at one end of a bundle of optical fibers. The light is transmitted through the bundle, and is finally magnified by the eyepiece.



m / Endoscopic images of a duodenal ulcer.

that your calculator will flash an error message at you when you try to take the inverse sine. What can happen physically in such a situation? The answer is that all the light is reflected, so there is no refracted ray. This phenomenon is known as *total internal reflection*, and is used in the fiber-optic cables that nowadays carry almost all long-distance telephone calls. The electrical signals from your phone travel to a switching center, where they are converted from electricity into light. From there, the light is sent across the country in a thin transparent fiber. The light is aimed straight into the end of the fiber, and as long as the fiber never goes through any turns that are too sharp, the light will always encounter the edge of the fiber at an angle sufficiently oblique to give total internal reflection. If the fiber-optic cable is thick enough, one can see an image at one end of whatever the other end is pointed at.

Alternatively, a bundle of cables can be used, since a single thick cable is too hard to bend. This technique for seeing around corners is useful for making surgery less traumatic. Instead of cutting a person wide open, a surgeon can make a small “keyhole” incision and insert a bundle of fiber-optic cable (known as an endoscope) into the body.

Since rays at sufficiently large angles with respect to the normal may be completely reflected, it is not surprising that the relative amount of reflection changes depending on the angle of incidence, and is greatest for large angles of incidence.

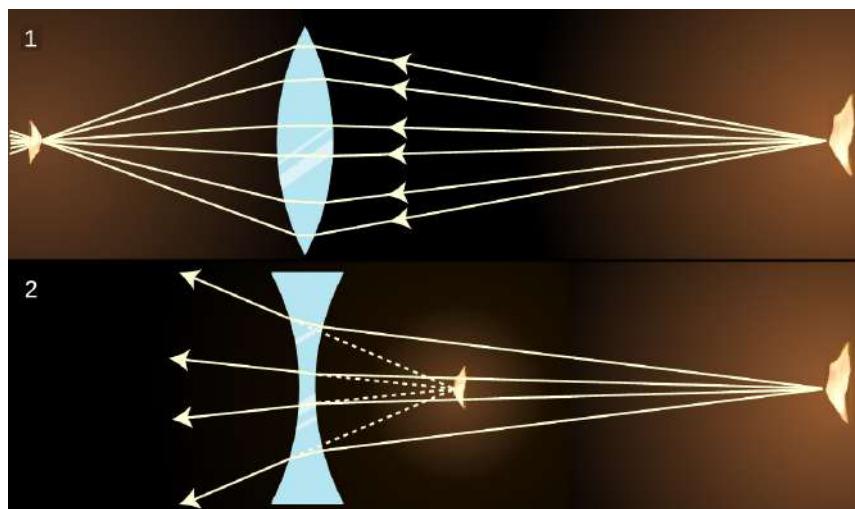
Discussion Questions

- A** What index of refraction should a fish have in order to be invisible to other fish?
- B** Does a surgeon using an endoscope need a source of light inside the body cavity? If so, how could this be done without inserting a light bulb through the incision?
- C** A denser sample of a gas has a higher index of refraction than a less dense sample (i.e., a sample under lower pressure), but why would it not make sense for the index of refraction of a gas to be proportional to density?
- D** The earth's atmosphere gets thinner and thinner as you go higher in altitude. If a ray of light comes from a star that is below the zenith, what will happen to it as it comes into the earth's atmosphere?
- E** Does total internal reflection occur when light in a denser medium encounters a less dense medium, or the other way around? Or can it occur in either case?

12.4.2 Lenses

Figures n/1 and n/2 show examples of lenses forming images. There is essentially nothing for you to learn about imaging with lenses that is truly new. You already know how to construct and use ray diagrams, and you know about real and virtual images. The

concept of the focal length of a lens is the same as for a curved mirror. The equations for locating images and determining magnifications are of the same form. It's really just a question of flexing your mental muscles on a few examples. The following self-checks and discussion questions will get you started. I've also made a video that demonstrates some applications and how to explain them with ray diagrams: <https://youtu.be/gL8awy6PWLQ>.



n / 1. A converging lens forms an image of a candle flame. 2. A diverging lens.

self-check F

- (1) In figures n/1 and n/2, classify the images as real or virtual.
- (2) Glass has an index of refraction that is greater than that of air. Consider the topmost ray in figure n/1. Explain why the ray makes a slight left turn upon entering the lens, and another left turn when it exits.
- (3) If the flame in figure n/2 was moved closer to the lens, what would happen to the location of the image? ▷ Answer, p. 1066

Discussion Questions

- A** In figures n/1 and n/2, the front and back surfaces are parallel to each other at the center of the lens. What will happen to a ray that enters near the center, but not necessarily along the axis of the lens? Draw a BIG ray diagram, and show a ray that comes from off axis.

In discussion questions B-F, don't draw ultra-detailed ray diagrams as in A.

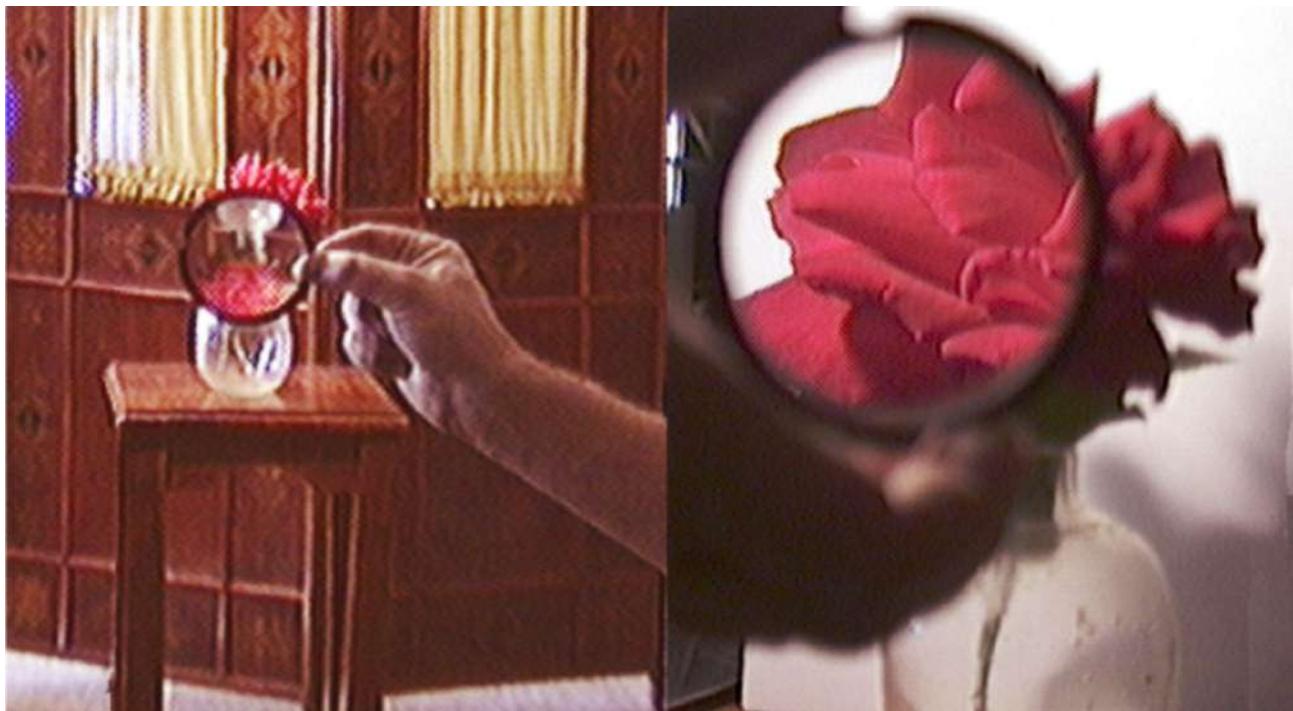
- B** Suppose you wanted to change the setup in figure n/1 so that the location of the actual flame in the figure would instead be occupied by an image of a flame. Where would you have to move the candle to achieve this? What about in n/2?

- C** There are three qualitatively different types of image formation that can occur with lenses, of which figures n/1 and n/2 exhaust only two. Figure out what the third possibility is. Which of the three possibilities can result in a magnification greater than one? Cf. problem 10, p. 831.

D Classify the examples shown in figure o according to the types of images delineated in discussion question C.

E In figures n/1 and n/2, the only rays drawn were those that happened to enter the lenses. Discuss this in relation to figure o.

F In the right-hand side of figure o, the image viewed through the lens is in focus, but the side of the rose that sticks out from behind the lens is not. Why?



o / Two images of a rose created by the same lens and recorded with the same camera.

12.4.3 ★ The lensmaker's equation

The focal length of a spherical mirror is simply $r/2$, but we cannot expect the focal length of a lens to be given by pure geometry, since it also depends on the index of refraction of the lens. Suppose we have a lens whose front and back surfaces are both spherical. (This is no great loss of generality, since any surface with a sufficiently shallow curvature can be approximated with a sphere.) Then if the lens is immersed in a medium with an index of refraction of 1, its focal length is given approximately by

$$f = \left[(n - 1) \left| \frac{1}{r_1} \pm \frac{1}{r_2} \right| \right]^{-1},$$

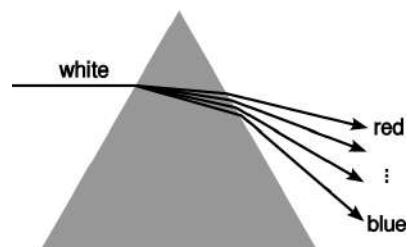
where n is the index of refraction and r_1 and r_2 are the radii of curvature of the two surfaces of the lens. This is known as the lensmaker's equation. In my opinion it is not particularly worthy

p / The radii of curvature appearing in the lensmaker's equation.

of memorization. The positive sign is used when both surfaces are curved outward or both are curved inward; otherwise a negative sign applies. The proof of this equation is left as an exercise to those readers who are sufficiently brave and motivated.

12.4.4 Dispersion

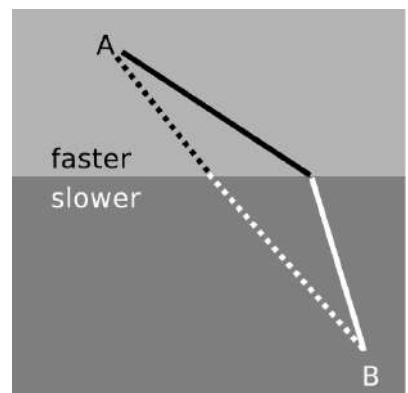
For most materials, we observe that the index of refraction depends slightly on wavelength, being highest at the blue end of the visible spectrum and lowest at the red. For example, white light disperses into a rainbow when it passes through a prism, q. Even when the waves involved aren't light waves, and even when refraction isn't of interest, the dependence of wave speed on wavelength is referred to as dispersion. Dispersion inside spherical raindrops is responsible for the creation of rainbows in the sky, and in an optical instrument such as the eye or a camera it is responsible for a type of aberration called chromatic aberration (subsection 12.3.3 and problem 28). As we'll see in subsection 13.3.2, dispersion causes a wave that is not a pure sine wave to have its shape distorted as it travels, and also causes the speed at which energy and information are transported by the wave to be different from what one might expect from a naive calculation. The microscopic reasons for dispersion of light in matter are discussed in optional subsection 12.4.6.



q / Dispersion of white light by a prism. White light is a mixture of all the wavelengths of the visible spectrum. Waves of different wavelengths undergo different amounts of refraction.

12.4.5 * The principle of least time for refraction

We have seen previously how the rules governing straight-line motion of light and reflection of light can be derived from the principle of least time. What about refraction? In the figure, it is indeed plausible that the bending of the ray serves to minimize the time required to get from a point A to point B. If the ray followed the unbent path shown with a dashed line, it would have to travel a longer distance in the medium in which its speed is slower. By bending the correct amount, it can reduce the distance it has to cover in the slower medium without going too far out of its way. It is true that Snell's law gives exactly the set of angles that minimizes the time required for light to get from one point to another. The proof of this fact is left as an exercise (problem 38, p. 836).

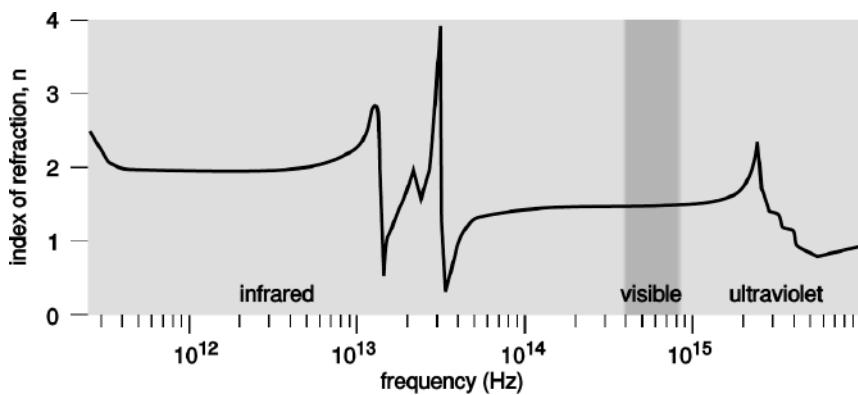


r / The principle of least time applied to refraction.

12.4.6 ★ Microscopic description of refraction

Given that the speed of light is different in different media, we've seen two different explanations (on p. 806 and in subsection 12.4.5 above) of why refraction must occur. What we haven't yet explained is why the speed of light does depend on the medium.

s / Index of refraction of silica glass, redrawn from Kitamura, Pilon, and Jonasz, Applied Optics 46 (2007) 8118, reprinted online at <http://www.seas.ucla.edu/~pilon/Publications/AO2007-1.pdf>.



A good clue as to what's going on comes from the figure s. The relatively minor variation of the index of refraction within the visible spectrum was misleading. At certain specific frequencies, n exhibits wild swings in the positive and negative directions. After each such swing, we reach a new, lower plateau on the graph. These frequencies are resonances. For example, the visible part of the spectrum lies on the left-hand tail of a resonance at about 2×10^{15} Hz, corresponding to the ultraviolet part of the spectrum. This resonance arises from the vibration of the electrons, which are bound to the nuclei as if by little springs. Because this resonance is narrow, the effect on visible-light frequencies is relatively small, but it is stronger at the blue end of the spectrum than at the red end. Near each resonance, not only does the index of refraction fluctuate wildly, but the glass becomes nearly opaque; this is because the vibration becomes very strong, causing energy to be dissipated as heat. The "staircase" effect is the same one visible in any resonance, e.g., figure k on p. 184: oscillators have a finite response for $f \ll f_0$, but the response approaches zero for $f \gg f_0$.

So far, we have a qualitative explanation of the frequency-variation of the loosely defined "strength" of the glass's effect on a light wave, but we haven't explained why the effect is observed as a change in speed, or why each resonance is an up-down swing rather than a single positive peak. To understand these effects in more detail, we need to consider the phase response of the oscillator. As shown in the bottom panel of figure j on p. 185, the phase response reverses itself as we pass through a resonance.

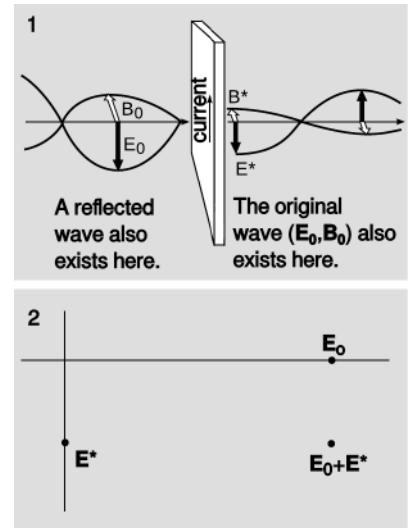
Suppose that a plane wave is normally incident on the left side of a thin sheet of glass, $t/1$, at $f \ll f_0$. The light wave observed on the

right side consists of a superposition of the incident wave consisting of \mathbf{E}_0 and \mathbf{B}_0 with a secondary wave \mathbf{E}^* and \mathbf{B}^* generated by the oscillating charges in the glass. Since the frequency is far below resonance, the response $q\mathbf{x}$ of a vibrating charge q is in phase with the driving force \mathbf{E}_0 . The current is the derivative of this quantity, and therefore 90 degrees ahead of it in phase. The magnetic field generated by a sheet of current has been analyzed in subsection 11.2.1, and the result, shown in figure e on p. 692, is just what we would expect from the right-hand rule. We find, t/1, that the secondary wave is 90 degrees ahead of the incident one in phase. The incident wave still exists on the right side of the sheet, but it is superposed with the secondary one. Their addition is shown in t/2 using the complex number representation introduced in subsection 10.5.7. The superposition of the two fields lags behind the incident wave, which is the effect we would expect if the wave had traveled more slowly through the glass.

In the case $f \gg f_0$, the same analysis applies except that the phase of the secondary wave is reversed. The transmitted wave is advanced rather than retarded in phase. This explains the dip observed in figure s after each spike.

All of this is in accord with our understanding of relativity, ch. 7, in which we saw that the universal speed c was to be understood fundamentally as a conversion factor between the units used to measure time and space — not as the speed of light. Since c isn't defined as the speed of light, it's of no fundamental importance whether light has a different speed in matter than it does in vacuum. In fact, the picture we've built up here is one in which all of our electromagnetic waves travel at c ; propagation at some other speed is only what appears to happen because of the superposition of the $(\mathbf{E}_0, \mathbf{B}_0)$ and $(\mathbf{E}^*, \mathbf{B}^*)$ waves, both of which move at c .

But it is worrisome that at the frequencies where $n < 1$, the speed of the wave is greater than c . According to special relativity, information is never supposed to be transmitted at speeds greater than c , since this would produce situations in which a signal could be received before it was transmitted! This difficulty is resolved in subsection 13.3.2, where we show that there are two different velocities that can be defined for a wave in a dispersive medium, the phase velocity and the group velocity. The group velocity is the velocity at which information is transmitted, and it is always less than c .



t / 1. A wave incident on a sheet of glass excites current in the glass, which produce a secondary wave. 2. The secondary wave superposes with the original wave, as represented in the complex-number representation introduced in subsection 10.5.7.

12.5 Wave optics

Electron microscopes can make images of individual atoms, but why will a visible-light microscope never be able to? Stereo speakers create the illusion of music that comes from a band arranged in your living room, but why doesn't the stereo illusion work with bass notes? Why are computer chip manufacturers investing billions of dollars in equipment to etch chips with x-rays instead of visible light?

The answers to all of these questions have to do with the subject of wave optics. So far this book has discussed the interaction of light waves with matter, and its practical applications to optical devices like mirrors, but we have used the ray model of light almost exclusively. Hardly ever have we explicitly made use of the fact that light is an electromagnetic wave. We were able to get away with the simple ray model because the chunks of matter we were discussing, such as lenses and mirrors, were thousands of times larger than a wavelength of light. We now turn to phenomena and devices that can only be understood using the wave model of light.

12.5.1 Diffraction

Figure a shows a typical problem in wave optics, enacted with water waves. It may seem surprising that we don't get a simple pattern like figure b, but the pattern would only be that simple if the wavelength was hundreds of times shorter than the distance between the gaps in the barrier and the widths of the gaps.

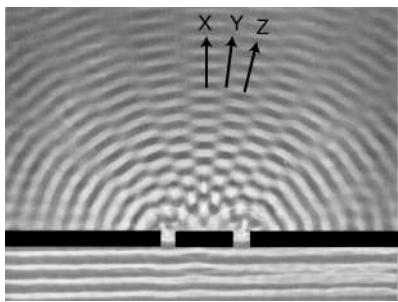
Wave optics is a broad subject, but this example will help us to pick out a reasonable set of restrictions to make things more manageable:

(1) We restrict ourselves to cases in which a wave travels through a uniform medium, encounters a certain area in which the medium has different properties, and then emerges on the other side into a second uniform region.

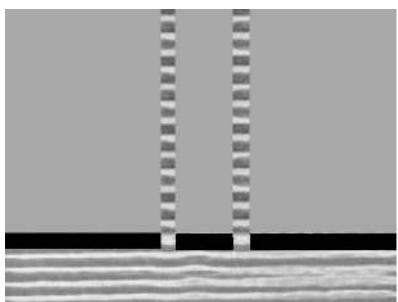
(2) We assume that the incoming wave is a nice tidy sine-wave pattern with wavefronts that are lines (or, in three dimensions, planes).

(3) In figure a we can see that the wave pattern immediately beyond the barrier is rather complex, but farther on it sorts itself out into a set of wedges separated by gaps in which the water is still. We will restrict ourselves to studying the simpler wave patterns that occur farther away, so that the main question of interest is how intense the outgoing wave is at a given angle.

The kind of phenomenon described by restriction (1) is called *diffraction*. Diffraction can be defined as the behavior of a wave when it encounters an obstacle or a nonuniformity in its medium. In general, diffraction causes a wave to bend around obstacles and



a / In this view from overhead, a straight, sinusoidal water wave encounters a barrier with two gaps in it. Strong wave vibration occurs at angles X and Z, but there is none at all at angle Y. (The figure has been retouched from a real photo of water waves. In reality, the waves beyond the barrier would be much weaker than the ones before it, and they would therefore be difficult to see.)



b / This doesn't happen.

make patterns of strong and weak waves radiating out beyond the obstacle. Understanding diffraction is the central problem of wave optics. If you understand diffraction, even the subset of diffraction problems that fall within restrictions (2) and (3), the rest of wave optics is icing on the cake.

Diffraction can be used to find the structure of an unknown diffracting object: even if the object is too small to study with ordinary imaging, it may be possible to work backward from the diffraction pattern to learn about the object. The structure of a crystal, for example, can be determined from its x-ray diffraction pattern.

Diffraction can also be a bad thing. In a telescope, for example, light waves are diffracted by all the parts of the instrument. This will cause the image of a star to appear fuzzy even when the focus has been adjusted correctly. By understanding diffraction, one can learn how a telescope must be designed in order to reduce this problem — essentially, it should have the biggest possible diameter.

There are two ways in which restriction (2) might commonly be violated. First, the light might be a mixture of wavelengths. If we simply want to observe a diffraction pattern or to use diffraction as a technique for studying the object doing the diffracting (e.g., if the object is too small to see with a microscope), then we can pass the light through a colored filter before diffracting it.

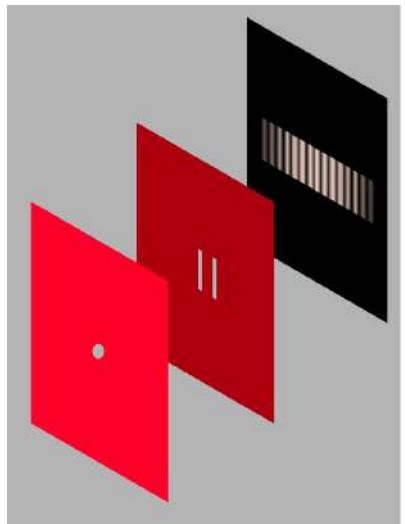
A second issue is that light from sources such as the sun or a lightbulb does not consist of a nice neat plane wave, except over very small regions of space. Different parts of the wave are out of step with each other, and the wave is referred to as *incoherent*. One way of dealing with this is shown in figure c. After filtering to select a certain wavelength of red light, we pass the light through a small pinhole. The region of the light that is intercepted by the pinhole is so small that one part of it is not out of step with another. Beyond the pinhole, light spreads out in a spherical wave; this is analogous to what happens when you speak into one end of a paper towel roll and the sound waves spread out in all directions from the other end. By the time the spherical wave gets to the double slit it has spread out and reduced its curvature, so that we can now think of it as a simple plane wave.

If this seems laborious, you may be relieved to know that modern technology gives us an easier way to produce a single-wavelength, coherent beam of light: the laser.

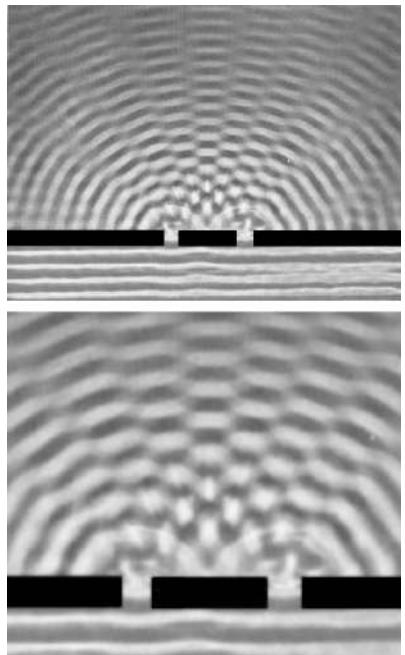
The parts of the final image on the screen in c are called diffraction fringes. The center of each fringe is a point of maximum brightness, and halfway between two fringes is a minimum.

Discussion Question

A Why would x-rays rather than visible light be used to find the structure



c / A practical, low-tech setup for observing diffraction of light.



d / The bottom figure is simply a copy of the middle portion of the top one, scaled up by a factor of two. All the angles are the same. Physically, the angular pattern of the diffraction fringes can't be any different if we scale both λ and d by the same factor, leaving λ/d unchanged.

of a crystal? Sound waves are used to make images of fetuses in the womb. What would influence the choice of wavelength?

12.5.2 Scaling of diffraction

This chapter has “optics” in its title, so it is nominally about light, but we started out with an example involving water waves. Water waves are certainly easier to visualize, but is this a legitimate comparison? In fact the analogy works quite well, despite the fact that a light wave has a wavelength about a million times shorter. This is because diffraction effects scale uniformly. That is, if we enlarge or reduce the whole diffraction situation by the same factor, including both the wavelengths and the sizes of the obstacles the wave encounters, the result is still a valid solution.

This is unusually simple behavior! In subsection 0.2.2 we saw many examples of more complex scaling, such as the impossibility of bacteria the size of dogs, or the need for an elephant to eliminate heat through its ears because of its small surface-to-volume ratio, whereas a tiny shrew’s life-style centers around conserving its body heat.

Of course water waves and light waves differ in many ways, not just in scale, but the general facts you will learn about diffraction are applicable to all waves. In some ways it might have been more appropriate to insert this chapter after section 6.2 on bounded waves, but many of the important applications are to light waves, and you would probably have found these much more difficult without any background in optics.

A portrait painting of Christiaan Huygens, a Dutch mathematician, physicist, and astronomer. He is shown from the chest up, wearing a white cravat and a dark jacket. He has dark hair and is looking slightly to the right.

Another way of stating the simple scaling behavior of diffraction is that the diffraction angles we get depend only on the unitless ratio λ/d , where λ is the wavelength of the wave and d is some dimension of the diffracting objects, e.g., the center-to-center spacing between the slits in figure a. If, for instance, we scale up both λ and d by a factor of 37, the ratio λ/d will be unchanged.

12.5.3 The correspondence principle

The only reason we don’t usually notice diffraction of light in everyday life is that we don’t normally deal with objects that are comparable in size to a wavelength of visible light, which is about a millionth of a meter. Does this mean that wave optics contradicts ray optics, or that wave optics sometimes gives wrong results? No. If you hold three fingers out in the sunlight and cast a shadow with them, *either* wave optics or ray optics can be used to predict the straightforward result: a shadow pattern with two bright lines where the light has gone through the gaps between your fingers. Wave optics is a more general theory than ray optics, so in any case where ray optics is valid, the two theories will agree. This is an example of a general idea enunciated by the physicist Niels Bohr, called the *correspondence principle*: when flaws in a physical theory

e / Christiaan Huygens (1629–1695).

lead to the creation of a new and more general theory, the new theory must still agree with the old theory within its more restricted area of applicability. After all, a theory is only created as a way of describing experimental observations. If the original theory had not worked in any cases at all, it would never have become accepted.

In the case of optics, the correspondence principle tells us that when λ/d is small, both the ray and the wave model of light must give approximately the same result. Suppose you spread your fingers and cast a shadow with them using a coherent light source. The quantity λ/d is about 10^{-4} , so the two models will agree very closely. (To be specific, the shadows of your fingers will be outlined by a series of light and dark fringes, but the angle subtended by a fringe will be on the order of 10^{-4} radians, so they will be too tiny to be visible.)

self-check G

What kind of wavelength would an electromagnetic wave have to have in order to diffract dramatically around your body? Does this contradict the correspondence principle?

▷ Answer, p. 1066

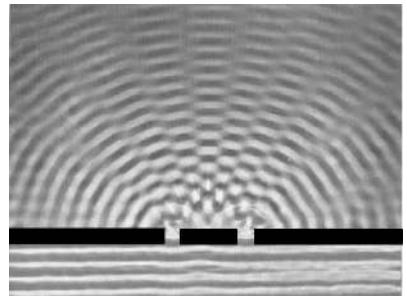
12.5.4 Huygens' principle

Returning to the example of double-slit diffraction, f, note the strong visual impression of two overlapping sets of concentric semi-circles. This is an example of *Huygens' principle*, named after a Dutch physicist and astronomer. (The first syllable rhymes with “boy.”) Huygens' principle states that any wavefront can be broken down into many small side-by-side wave peaks, g, which then spread out as circular ripples, h, and by the principle of superposition, the result of adding up these sets of ripples must give the same result as allowing the wave to propagate forward, i. In the case of sound or light waves, which propagate in three dimensions, the “ripples” are actually spherical rather than circular, but we can often imagine things in two dimensions for simplicity.

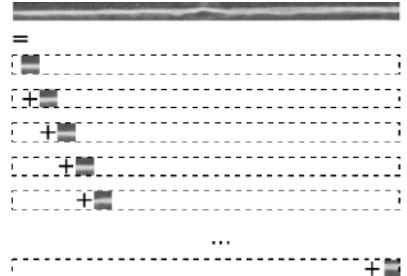
In double-slit diffraction the application of Huygens' principle is visually convincing: it is as though all the sets of ripples have been blocked except for two. It is a rather surprising mathematical fact, however, that Huygens' principle gives the right result in the case of an unobstructed linear wave, h and i. A theoretically infinite number of circular wave patterns somehow conspire to add together and produce the simple linear wave motion with which we are familiar.

Since Huygens' principle is equivalent to the principle of superposition, and superposition is a property of waves, what Huygens had created was essentially the first wave theory of light. However, he imagined light as a series of pulses, like hand claps, rather than as a sinusoidal wave.

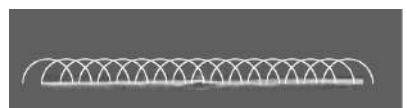
The history is interesting. Isaac Newton loved the atomic theory of matter so much that he searched enthusiastically for evidence that



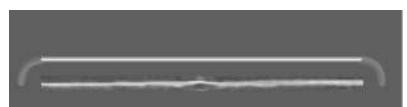
f / Double-slit diffraction.



g / A wavefront can be analyzed by the principle of superposition, breaking it down into many small parts.



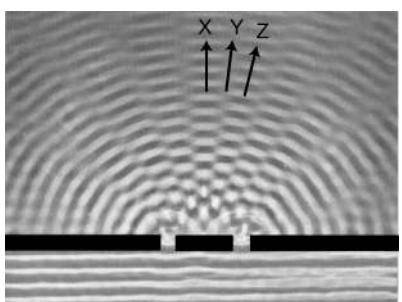
h / If it was by itself, each of the parts would spread out as a circular ripple.



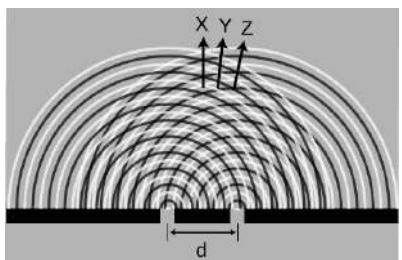
i / Adding up the ripples produces a new wavefront.



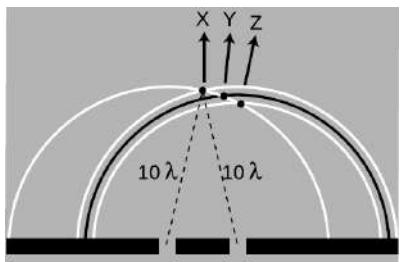
j / Thomas Young



k / Double-slit diffraction.



l / Use of Huygens' principle.



m / Constructive interference along the center-line.

light was also made of tiny particles. The paths of his light particles would correspond to rays in our description; the only significant difference between a ray model and a particle model of light would occur if one could isolate individual particles and show that light had a “graininess” to it. Newton never did this, so although he thought of his model as a particle model, it is more accurate to say he was one of the builders of the ray model.

Almost all that was known about reflection and refraction of light could be interpreted equally well in terms of a particle model or a wave model, but Newton had one reason for strongly opposing Huygens' wave theory. Newton knew that waves exhibited diffraction, but diffraction of light is difficult to observe, so Newton believed that light did not exhibit diffraction, and therefore must not be a wave. Although Newton's criticisms were fair enough, the debate also took on the overtones of a nationalistic dispute between England and continental Europe, fueled by English resentment over Leibniz's supposed plagiarism of Newton's calculus. Newton wrote a book on optics, and his prestige and political prominence tended to discourage questioning of his model.

Thomas Young (1773-1829) was the person who finally, a hundred years later, did a careful search for wave interference effects with light and analyzed the results correctly. He observed double-slit diffraction of light as well as a variety of other diffraction effects, all of which showed that light exhibited wave interference effects, and that the wavelengths of visible light waves were extremely short. The crowning achievement was the demonstration by the experimentalist Heinrich Hertz and the theorist James Clerk Maxwell that light was an *electromagnetic* wave. Maxwell is said to have related his discovery to his wife one starry evening and told her that she was the only other person in the world who knew what starlight was.

12.5.5 Double-slit diffraction

Let's now analyze double-slit diffraction, k, using Huygens' principle. The most interesting question is how to compute the angles such as X and Z where the wave intensity is at a maximum, and the in-between angles like Y where it is minimized. Let's measure all our angles with respect to the vertical center line of the figure, which was the original direction of propagation of the wave.

If we assume that the width of the slits is small (on the order of the wavelength of the wave or less), then we can imagine only a single set of Huygens ripples spreading out from each one, l. White lines represent peaks, black ones troughs. The only dimension of the diffracting slits that has any effect on the geometric pattern of the overlapping ripples is then the center-to-center distance, d , between the slits.

We know from our discussion of the scaling of diffraction that there must be some equation that relates an angle like θ_Z to the ratio λ/d ,

$$\frac{\lambda}{d} \leftrightarrow \theta_Z.$$

If the equation for θ_Z depended on some other expression such as $\lambda + d$ or λ^2/d , then it would change when we scaled λ and d by the same factor, which would violate what we know about the scaling of diffraction.

Along the central maximum line, X, we always have positive waves coinciding with positive ones and negative waves coinciding with negative ones. (I have arbitrarily chosen to take a snapshot of the pattern at a moment when the waves emerging from the slit are experiencing a positive peak.) The superposition of the two sets of ripples therefore results in a doubling of the wave amplitude along this line. There is constructive interference. This is easy to explain, because by symmetry, each wave has had to travel an equal number of wavelengths to get from its slit to the center line, m: Because both sets of ripples have ten wavelengths to cover in order to reach the point along direction X, they will be in step when they get there.

At the point along direction Y shown in the same figure, one wave has traveled ten wavelengths, and is therefore at a positive extreme, but the other has traveled only nine and a half wavelengths, so it is at a negative extreme. There is perfect cancellation, so points along this line experience no wave motion.

But the distance traveled does not have to be equal in order to get constructive interference. At the point along direction Z, one wave has gone nine wavelengths and the other ten. They are both at a positive extreme.

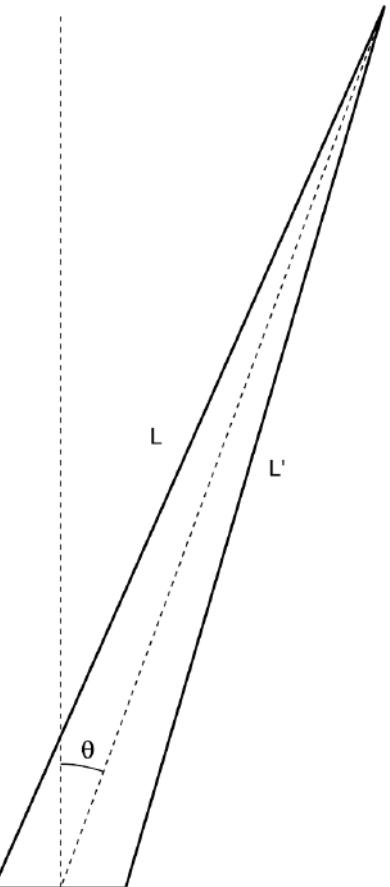
self-check H

At a point half a wavelength below the point marked along direction X, carry out a similar analysis. ▷ Answer, p. 1066

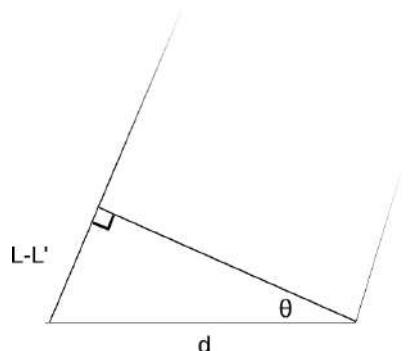
To summarize, we will have perfect constructive interference at any point where the distance to one slit differs from the distance to the other slit by an integer number of wavelengths. Perfect destructive interference will occur when the number of wavelengths of path length difference equals an integer plus a half.

Now we are ready to find the equation that predicts the angles of the maxima and minima. The waves travel different distances to get to the same point in space, n. We need to find whether the waves are in phase (in step) or out of phase at this point in order to predict whether there will be constructive interference, destructive interference, or something in between.

One of our basic assumptions in this chapter is that we will only be dealing with the diffracted wave in regions very far away from the



n / The waves travel distances L and L' from the two slits to get to the same point in space, at an angle θ from the center line.



o / A close-up view of figure n, showing how the path length difference $L - L'$ is related to d and to the angle θ .

object that diffracts it, so the triangle is long and skinny. Most real-world examples with diffraction of light, in fact, would have triangles with even skinnier proportions than this one. The two long sides are therefore very nearly parallel, and we are justified in drawing the right triangle shown in figure o, labeling one leg of the right triangle as the difference in path length, $L - L'$, and labeling the acute angle as θ . (In reality this angle is a tiny bit greater than the one labeled θ in figure n.)

The difference in path length is related to d and θ by the equation

$$\frac{L - L'}{d} = \sin \theta.$$

Constructive interference will result in a maximum at angles for which $L - L'$ is an integer number of wavelengths,

$$L - L' = m\lambda.$$

[condition for a maximum;
 m is an integer]

Here m equals 0 for the central maximum, -1 for the first maximum to its left, $+2$ for the second maximum on the right, etc. Putting all the ingredients together, we find $m\lambda/d = \sin \theta$, or

$$\frac{\lambda}{d} = \frac{\sin \theta}{m}.$$

[condition for a maximum;
 m is an integer]

Similarly, the condition for a minimum is

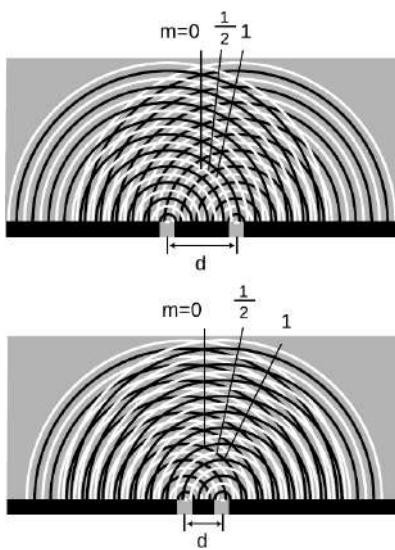
$$\frac{\lambda}{d} = \frac{\sin \theta}{m}.$$

[condition for a minimum;
 m is an integer plus 1/2]

That is, the minima are about halfway between the maxima.

As expected based on scaling, this equation relates angles to the unitless ratio λ/d . Alternatively, we could say that we have proven the scaling property in the special case of double-slit diffraction. It was inevitable that the result would have these scaling properties, since the whole proof was geometric, and would have been equally valid when enlarged or reduced on a photocopying machine!

Counterintuitively, this means that a diffracting object with smaller dimensions produces a bigger diffraction pattern, p.



p / Cutting d in half doubles the angles of the diffraction fringes.



q / Double-slit diffraction patterns of long-wavelength red light (top) and short-wavelength blue light (bottom).

Double-slit diffraction of blue and red light

example 12

Blue light has a shorter wavelength than red. For a given double-slit spacing d , the smaller value of λ/d for leads to smaller values of $\sin \theta$, and therefore to a more closely spaced set of diffraction fringes, as shown in figure q.

The correspondence principle

example 13

Let's also consider how the equations for double-slit diffraction relate to the correspondence principle. When the ratio λ/d is very small, we should recover the case of simple ray optics. Now if λ/d is small, $\sin \theta$ must be small as well, and the spacing between the diffraction fringes will be small as well. Although we have not proven it, the central fringe is always the brightest, and the fringes get dimmer and dimmer as we go farther from it. For small values of λ/d , the part of the diffraction pattern that is bright enough to be detectable covers only a small range of angles. This is exactly what we would expect from ray optics: the rays passing through the two slits would remain parallel, and would continue moving in the $\theta = 0$ direction. (In fact there would be images of the two separate slits on the screen, but our analysis was all in terms of angles, so we should not expect it to address the issue of whether there is structure within a set of rays that are all traveling in the $\theta = 0$ direction.)

Spacing of the fringes at small angles

example 14

At small angles, we can use the approximation $\sin \theta \approx \theta$, which is valid if θ is measured in radians. The equation for double-slit diffraction becomes simply

$$\frac{\lambda}{d} = \frac{\theta}{m},$$

which can be solved for θ to give

$$\theta = \frac{m\lambda}{d}.$$

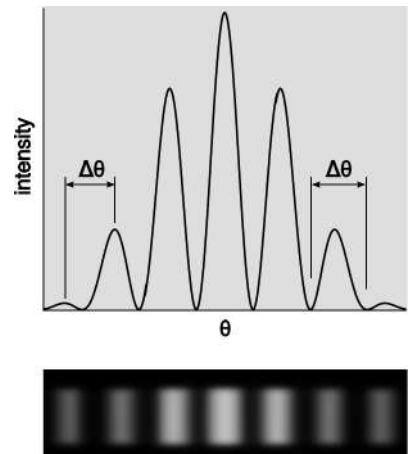
The difference in angle between successive fringes is the change in θ that results from changing m by plus or minus one,

$$\Delta\theta = \frac{\lambda}{d}.$$

For example, if we write θ_7 for the angle of the seventh bright fringe on one side of the central maximum and θ_8 for the neighboring one, we have

$$\begin{aligned}\theta_8 - \theta_7 &= \frac{8\lambda}{d} - \frac{7\lambda}{d} \\ &= \frac{\lambda}{d},\end{aligned}$$

and similarly for any other neighboring pair of fringes.



r / Interpretation of the angular spacing $\Delta\theta$ in example 14. It can be defined either from maximum to maximum or from minimum to minimum. Either way, the result is the same. It does not make sense to try to interpret $\Delta\theta$ as the width of a fringe; one can see from the graph and from the image below that it is not obvious either that such a thing is well defined or that it would be the same for all fringes.

Although the equation $\lambda/d = \sin \theta/m$ is only valid for a double slit, it is still a guide to our thinking even if we are observing diffraction of light by a virus or a flea's leg: it is always true that

- (1) large values of λ/d lead to a broad diffraction pattern, and
- (2) diffraction patterns are repetitive.

In many cases the equation looks just like $\lambda/d = \sin \theta/m$ but with an extra numerical factor thrown in, and with d interpreted as some other dimension of the object, e.g., the diameter of a piece of wire.

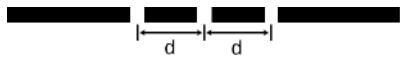
12.5.6 Repetition

Suppose we replace a double slit with a triple slit, s. We can think of this as a third repetition of the structures that were present in the double slit. Will this device be an improvement over the double slit for any practical reasons?

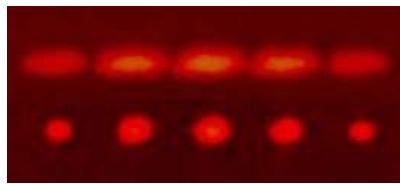
The answer is yes, as can be shown using figure u. For ease of visualization, I have violated our usual rule of only considering points very far from the diffracting object. The scale of the drawing is such that a wavelength is one cm. In u/1, all three waves travel an integer number of wavelengths to reach the same point, so there is a bright central spot, as we would expect from our experience with the double slit. In figure u/2, we show the path lengths to a new point. This point is farther from slit A by a quarter of a wavelength, and correspondingly closer to slit C. The distance from slit B has hardly changed at all. Because the paths traveled from slits A and C differ by half a wavelength, there will be perfect destructive interference between these two waves. There is still some uncanceled wave intensity because of slit B, but the amplitude will be three times less than in figure u/1, resulting in a factor of 9 decrease in brightness. Thus, by moving off to the right a little, we have gone from the bright central maximum to a point that is quite dark.

Now let's compare with what would have happened if slit C had been covered, creating a plain old double slit. The waves coming from slits A and B would have been out of phase by 0.23 wavelengths, but this would not have caused very severe interference. The point in figure u/2 would have been quite brightly lit up.

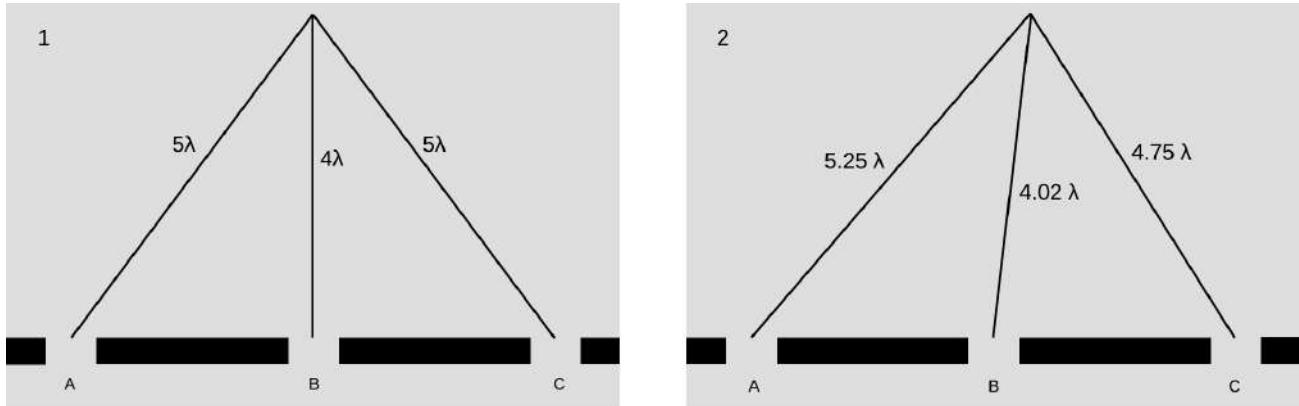
To summarize, we have found that adding a third slit narrows down the central fringe dramatically. The same is true for all the other fringes as well, and since the same amount of energy is con-



s / A triple slit.



t / A double-slit diffraction pattern (top), and a pattern made by five slits (bottom).



u / 1. There is a bright central maximum. 2. At this point just off the central maximum, the path lengths traveled by the three waves have changed.

centrated in narrower diffraction fringes, each fringe is brighter and easier to see, t.

This is an example of a more general fact about diffraction: if some feature of the diffracting object is repeated, the locations of the maxima and minima are unchanged, but they become narrower.

Taking this reasoning to its logical conclusion, a diffracting object with thousands of slits would produce extremely narrow fringes. Such an object is called a diffraction grating.

12.5.7 Single-slit diffraction

If we use only a single slit, is there diffraction? If the slit is not wide compared to a wavelength of light, then we can approximate its behavior by using only a single set of Huygens ripples. There are no other sets of ripples to add to it, so there are no constructive or destructive interference effects, and no maxima or minima. The result will be a uniform spherical wave of light spreading out in all directions, like what we would expect from a tiny lightbulb. We could call this a diffraction pattern, but it is a completely featureless one, and it could not be used, for instance, to determine the wavelength of the light, as other diffraction patterns could.

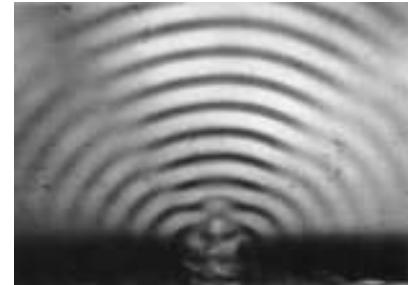
All of this, however, assumes that the slit is narrow compared to a wavelength of light. If, on the other hand, the slit is broader, there will indeed be interference among the sets of ripples spreading out from various points along the opening. Figure v shows an example with water waves, and figure w with light.

self-check 1

How does the wavelength of the waves compare with the width of the slit in figure v?

▷ Answer, p. 1066

We will not go into the details of the analysis of single-slit diffraction, but let us see how its properties can be related to the general



v / Single-slit diffraction of water waves.



w / Single-slit diffraction of red light. Note the double width of the central maximum.



x / A pretty good simulation of the single-slit pattern of figure v, made by using three motors to produce overlapping ripples from three neighboring points in the water.

things we've learned about diffraction. We know based on scaling arguments that the angular sizes of features in the diffraction pattern must be related to the wavelength and the width, a , of the slit by some relationship of the form

$$\frac{\lambda}{a} \leftrightarrow \theta.$$

This is indeed true, and for instance the angle between the maximum of the central fringe and the maximum of the next fringe on one side equals $1.5\lambda/a$. Scaling arguments will never produce factors such as the 1.5, but they tell us that the answer must involve λ/a , so all the familiar qualitative facts are true. For instance, shorter-wavelength light will produce a more closely spaced diffraction pattern.

An important scientific example of single-slit diffraction is in telescopes. Images of individual stars, as in figure y, are a good way to examine diffraction effects, because all stars except the sun are so far away that no telescope, even at the highest magnification, can image their disks or surface features. Thus any features of a star's image must be due purely to optical effects such as diffraction. A prominent cross appears around the brightest star, and dimmer ones surround the dimmer stars. Something like this is seen in most telescope photos, and indicates that inside the tube of the telescope there were two perpendicular struts or supports. Light diffracted around these struts. You might think that diffraction could be eliminated entirely by getting rid of all obstructions in the tube, but the circles around the stars are diffraction effects arising from single-slit diffraction at the mouth of the telescope's tube! (Actually we have not even talked about diffraction through a circular opening, but the idea is the same.) Since the angular sizes of the diffracted images depend on λ/a , the only way to improve the resolution of the images is to increase the diameter, a , of the tube. This is one of the main reasons (in addition to light-gathering power) why the best telescopes must be very large in diameter.

self-check J

What would this imply about radio telescopes as compared with visible-light telescopes?

▷ Answer, p.

1066

Double-slit diffraction is easier to understand conceptually than single-slit diffraction, but if you do a double-slit diffraction experiment in real life, you are likely to encounter a complicated pattern like figure aa/1, rather than the simpler one, 2, you were expecting. This is because the slits are fairly big compared to the wavelength of the light being used. We really have two different distances in our pair of slits: d , the distance between the slits, and w , the width of each slit. Remember that smaller distances on the object the light diffracts around correspond to larger features of the diffraction pattern. The pattern 1 thus has two spacings in it: a short spac-

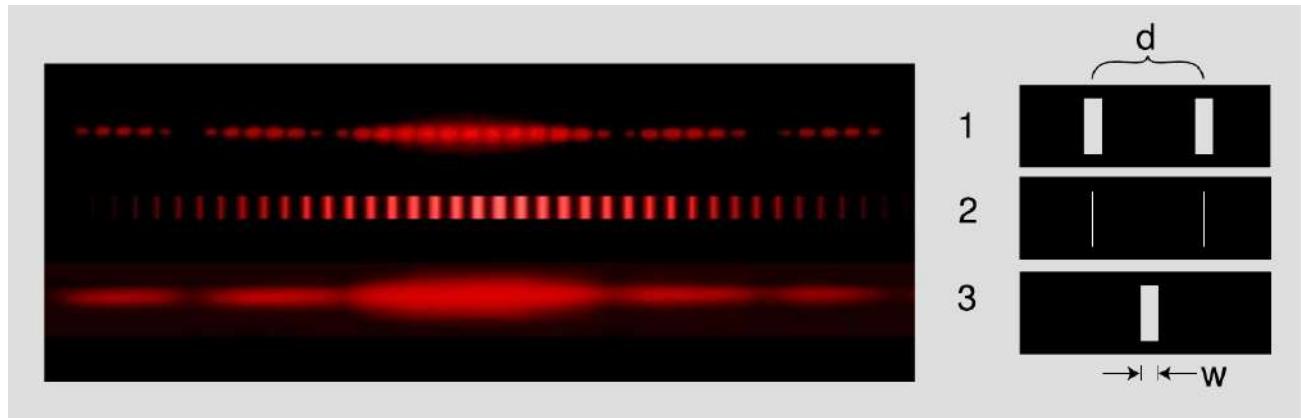


y / An image of the Pleiades star cluster. The circular rings around the bright stars are due to single-slit diffraction at the mouth of the telescope's tube.



z / A radio telescope.

ing corresponding to the large distance d , and a long spacing that relates to the small dimension w .



aa / 1. A diffraction pattern formed by a real double slit. The width of each slit is fairly big compared to the wavelength of the light. This is a real photo. 2. This idealized pattern is not likely to occur in real life. To get it, you would need each slit to be so narrow that its width was comparable to the wavelength of the light, but that's not usually possible. This is not a real photo. 3. A real photo of a single-slit diffraction pattern caused by a slit whose width is the same as the widths of the slits used to make the top pattern.

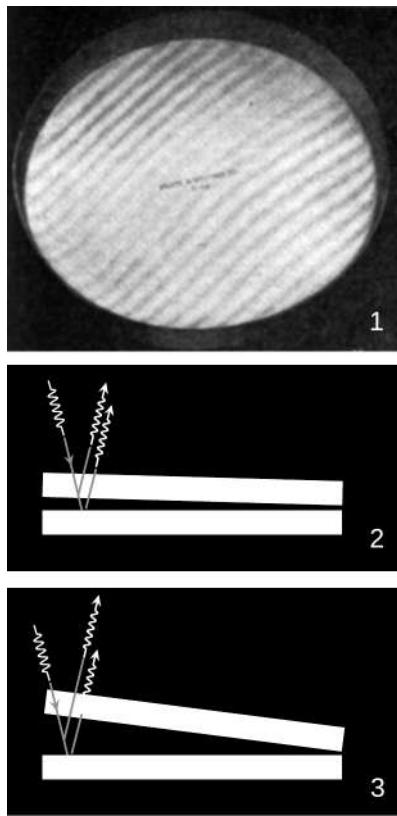
Discussion Question

A Why is it optically impossible for bacteria to evolve eyes that use visible light to form images?

12.5.8 Coherence

Up until now, we have avoided too much detailed discussion of two facts that sometimes make interference and diffraction effects unobservable, and that historically made them more difficult to discover. First there is the fact that white light is a mixture of all the visible wavelengths. This is why, for example, the thin-film interference pattern of a soap bubble looks like a rainbow. To simplify things, we need a source of light that is monochromatic, i.e., contains only a single wavelength or a small range of wavelengths. We could do this either by filtering a white light source or by using a source of light that is intrinsically monochromatic, such as a laser or some gas discharge tubes.

But even with a monochromatic light source, we encounter a separate issue, which is that most light sources do not emit light waves that are perfect, infinitely long sine waves. Sunlight and candlelight, for example, can be thought of as being composed of separate little spurts of light, referred to as wave packets or wave trains. Each wave packet is emitted by a separate atom of the gas. It contains some number of wavelengths, and it has no fixed phase relationship to any other wave packet. The wave trains emitted by a laser are much longer, but still not infinitely long.



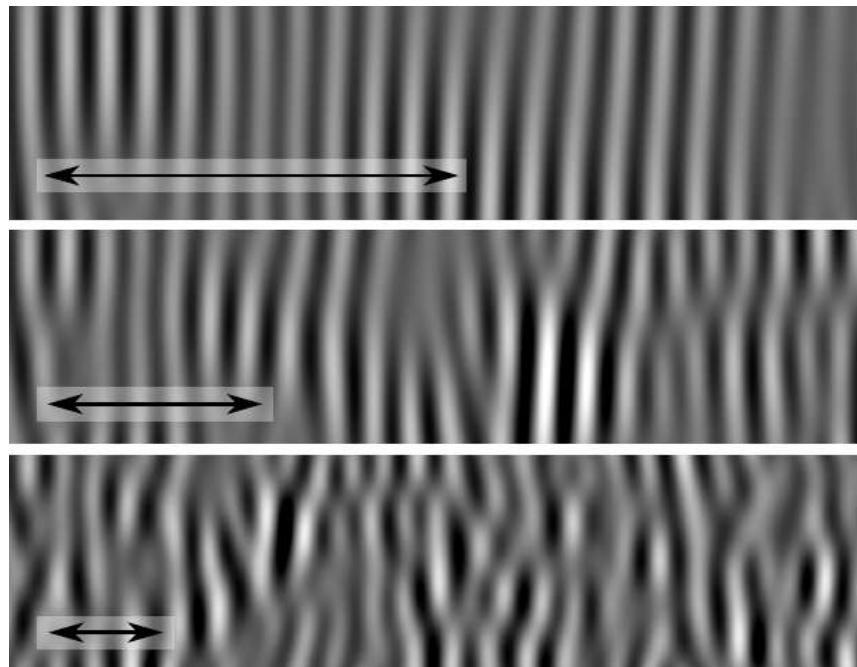
ab / 1. Interference in an air wedge. 2. Side view. 3. If the wedge is thicker than the coherence length of the light, the interference pattern disappears.

As an example of an experiment that can show these effects, figure ab/1 shows a thin-film interference pattern created by the air wedge between two pieces of very flat glass, where the top piece is placed at a very small angle relative to the bottom one, ab/2. The phase relationship between the two reflected waves is determined by the extra distance traveled by the ray that is reflected by the bottom plate (as well as the fact that one of the two reflections will be inverting).

If the angle is opened up too much, ab/3, we will no longer see fringes where the air layer is too thick. This is because the incident wave train has only a certain length, and the extra distance traveled is now so great that the two reflected wave trains no longer overlap in space. In general, if the incident wave trains are n wavelengths long, then we can see at most n bright and n dark fringes. The fact that about 18 fringes are visible in ab/1 shows that the light source used (let's say a sodium gas discharge tube) made wave trains at least 18 wavelengths in length.

In real-world light sources, the wave packets may not be as neat and tidy as the ones in figure ab. They may not look like sine waves with clean cut-offs at the ends, and they may overlap one another. The result will look more like the examples in figure ac. Such a wave pattern has a property called its coherence length L . On scales small compared to L , the wave appears like a perfect sine wave. On scales large compared to L , we lose all phase correlations. For example, the middle wave in figure ac has $L \approx 5\lambda$. If we pick two points within this wave separated by a distance of λ in the left-right direction, they are likely to be very nearly in phase. But if the

ac / Waves with three different coherence lengths, indicated by the arrows. Note that although there is a superficial similarity between these pictures and figure ab/1, they represent completely different things. Figure ab/1 is an actual photograph of interference fringes, whose brightness is proportional to the square of the amplitude. This figure is a picture of the wave's amplitude, not the squared amplitude, and is analogous to the little sine waves in ab/2 and ab/3. These are waves that are *traveling* across the page at the speed of light.



separation is 20λ , approximately the width of the entire figure, the phase relationship is essentially random. If the light comes from a flame or a gas discharge tube, then this lack of a phase relationship would be because the parts of the wave at these large separations from one another probably originated from different atoms in the source.

12.5.9 * The principle of least time

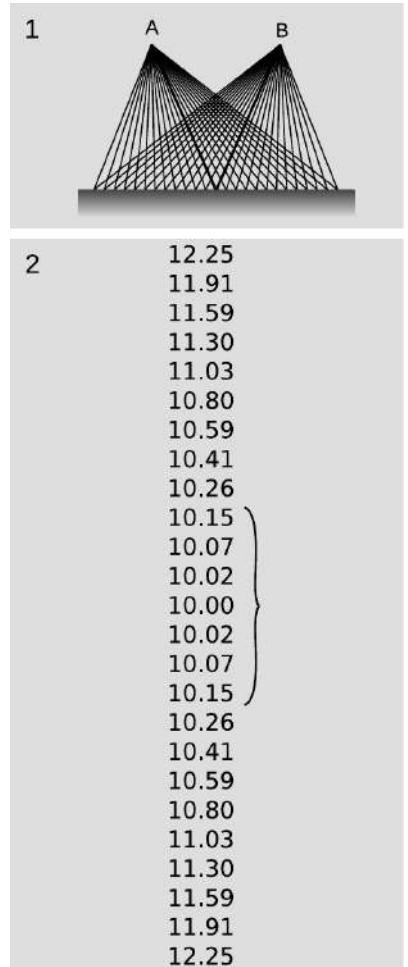
In subsection 12.1.5 and 12.4.5, we saw how in the ray model of light, both refraction and reflection can be described in an elegant and beautiful way by a single principle, the principle of least time. We can now justify the principle of least time based on the wave model of light. Consider an example involving reflection, ad. Starting at point A, Huygens' principle for waves tells us that we can think of the wave as spreading out in all directions. Suppose we imagine all the possible ways that a ray could travel from A to B. We show this by drawing 25 possible paths, of which the central one is the shortest. Since the principle of least time connects the wave model to the ray model, we should expect to get the most accurate results when the wavelength is much shorter than the distances involved — for the sake of this numerical example, let's say that a wavelength is $1/10$ of the shortest reflected path from A to B. The table, 2, shows the distances traveled by the 25 rays.

Note how similar are the distances traveled by the group of 7 rays, indicated with a bracket, that come closest to obeying the principle of least time. If we think of each one as a wave, then all 7 are again nearly in phase at point B. However, the rays that are farther from satisfying the principle of least time show more rapidly changing distances; on reuniting at point B, their phases are a random jumble, and they will very nearly cancel each other out. Thus, almost none of the wave energy delivered to point B goes by these longer paths. Physically we find, for instance, that a wave pulse emitted at A is observed at B after a time interval corresponding very nearly to the shortest possible path, and the pulse is not very “smeared out” when it gets there. The shorter the wavelength compared to the dimensions of the figure, the more accurate these approximate statements become.

Instead of drawing a finite number of rays, such as 25, what happens if we think of the angle, θ , of emission of the ray as a continuously varying variable? Minimizing the distance L requires

$$\frac{dL}{d\theta} = 0.$$

Because L is changing slowly in the vicinity of the angle that satisfies the principle of least time, all the rays that come out close to this angle have very nearly the same L , and remain very nearly in phase when they reach B. This is the basic reason why the discrete



ad / Light could take many different paths from A to B.

table, $ad/2$, turned out to have a group of rays that all traveled nearly the same distance.

As discussed in subsection 12.1.5, the principle of least time is really a principle of least *or greatest* time. This makes perfect sense, since $dL/d\theta = 0$ can in general describe either a minimum or a maximum

The principle of least time is very general. It does not apply just to refraction and reflection — it can even be used to prove that light rays travel in a straight line through empty space, without taking detours! This general approach to wave motion was used by Richard Feynman, one of the pioneers who in the 1950's reconciled quantum mechanics with relativity. A very readable explanation is given in a book Feynman wrote for laypeople, QED: The Strange Theory of Light and Matter.

Problems

The symbols \checkmark , \blacksquare , etc. are explained on page 846.

1 Draw a ray diagram showing why a small light source (a candle, say) produces sharper shadows than a large one (e.g., a long fluorescent bulb). \blacksquare

2 A Global Positioning System (GPS) receiver is a device that lets you figure out where you are by receiving timed radio signals from satellites. It works by measuring the travel time for the signals, which is related to the distance between you and the satellite. By finding the ranges to several different satellites in this way, it can pin down your location in three dimensions to within a few meters. How accurate does the measurement of the time delay have to be to determine your position to this accuracy? \blacksquare

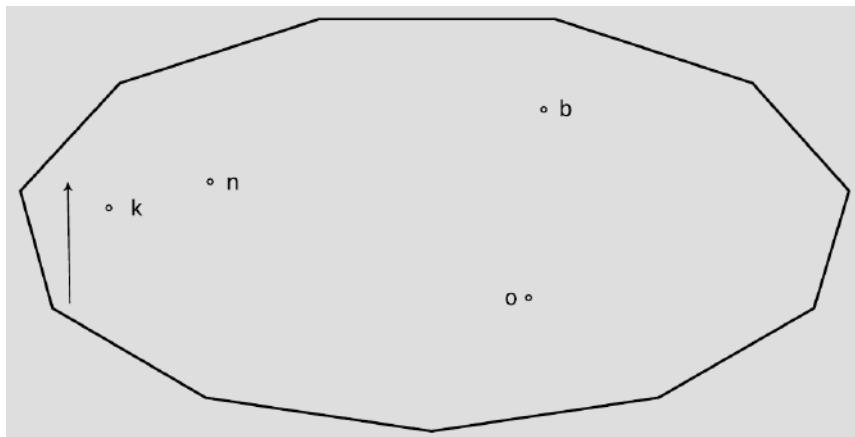
3 Estimate the frequency of an electromagnetic wave whose wavelength is similar in size to an atom (about a nm). Referring back to figure 0 on p. 732, in what part of the electromagnetic spectrum would such a wave lie (infrared, gamma-rays, . . .)? \blacksquare

4 The Stealth Bomber is designed with flat, smooth surfaces. Why would this make it difficult to detect using radar?

\triangleright Solution, p. 1049 \blacksquare

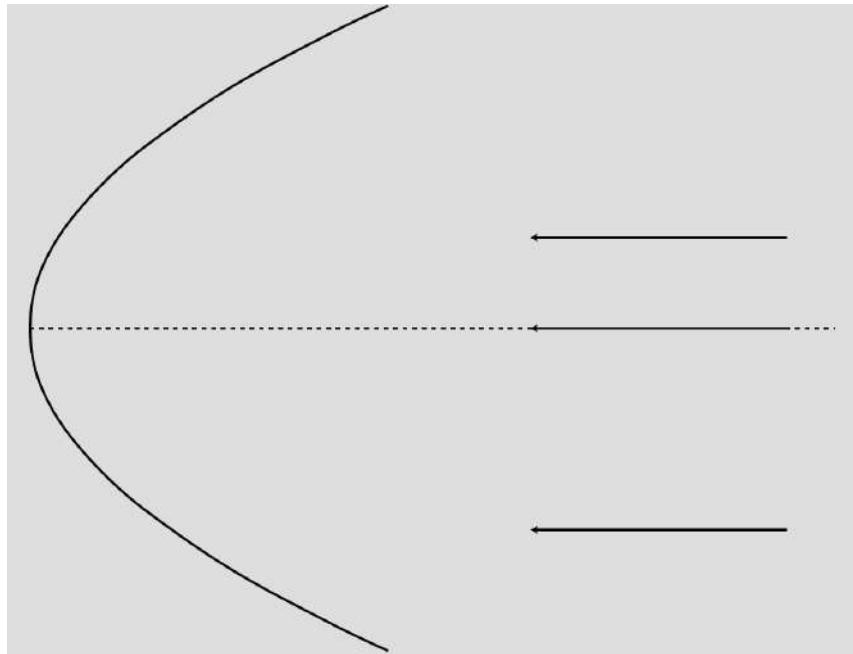
5 The natives of planet Wumpus play pool using light rays on an eleven-sided table with mirrors for bumpers, shown in the figure on the next page. Trace this shot accurately with a ruler to reveal the hidden message. To get good enough accuracy, you'll need to photocopy the page (or download the book and print the page) and construct each reflection using a protractor.

\triangleright Solution, p. 1050 \blacksquare



Problem 5.

6 The figure on the next page shows a curved (parabolic) mirror, with three parallel light rays coming toward it. One ray is approaching along the mirror's center line. (a) Continue the light rays until they are about to undergo their second reflection. To get good enough accuracy, you'll need to photocopy the page (or download the book and print the page) and draw in the normal at each place where a ray is reflected. What do you notice? (b) Make up an example of a practical use for this device. (c) How could you use this mirror with a small lightbulb to produce a parallel beam of light rays going off to the right? ▷ Solution, p. 1050 ■



Problem 6.

7 A man is walking at 1.0 m/s directly towards a flat mirror. At what speed is his separation from his image decreasing? ✓ ■

8 If a mirror on a wall is only big enough for you to see yourself from your head down to your waist, can you see your entire body by backing up? Test this experimentally and come up with an explanation for your observations, including a ray diagram.

Note that when you do the experiment, it's easy to confuse yourself if the mirror is even a tiny bit off of vertical. One way to check yourself is to artificially lower the top of the mirror by putting a piece of tape or a post-it note where it blocks your view of the top of your head. You can then check whether you are able to see more of yourself both above *and* below by backing up. ■

9 In section 12.2 we've only done examples of mirrors with hollowed-out shapes (called concave mirrors). Now draw a ray diagram for a curved mirror that has a bulging outward shape (called a convex mirror). (a) How does the image's distance from the mirror compare with the actual object's distance from the mirror? From this comparison, determine whether the magnification is greater than or less than one. (b) Is the image real, or virtual? Could this mirror ever make the other type of image? □

10 As discussed in question 9, there are two types of curved mirrors, concave and convex. Make a list of all the possible combinations of types of images (virtual or real) with types of mirrors (concave and convex). (Not all of the four combinations are physically possible.) Now for each one, use ray diagrams to determine whether increasing the distance of the object from the mirror leads to an increase or a decrease in the distance of the image from the mirror.

Draw BIG ray diagrams! Each diagram should use up about half a page of paper.

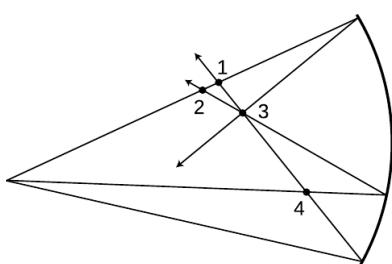
Some tips: To draw a ray diagram, you need two rays. For one of these, pick the ray that comes straight along the mirror's axis, since its reflection is easy to draw. After you draw the two rays and locate the image for the original object position, pick a new object position that results in the same type of image, and start a new ray diagram, in a different color of pen, right on top of the first one. For the two new rays, pick the ones that just happen to hit the mirror at the same two places; this makes it much easier to get the result right without depending on extreme accuracy in your ability to draw the reflected rays. □

11 If the user of an astronomical telescope moves her head closer to or farther away from the image she is looking at, does the magnification change? Does the angular magnification change? Explain. (For simplicity, assume that no eyepiece is being used.)

▷ Solution, p. 1050 □

12 In figure g/2 in on page 784, only the image of my forehead was located by drawing rays. Either photocopy the figure or download the book and print out the relevant page. On this copy of the figure, make a new set of rays coming from my chin, and locate its image. To make it easier to judge the angles accurately, draw rays from the chin that happen to hit the mirror at the same points where the two rays from the forehead were shown hitting it. By comparing the locations of the chin's image and the forehead's image, verify that the image is actually upside-down, as shown in the original figure. □

13 The figure shows four points where rays cross. Of these, which are image points? Explain. □



Problem 13.

14 Here's a game my kids like to play. I sit next to a sunny window, and the sun reflects from the glass on my watch, making a disk of light on the wall or floor, which they pretend to chase as I move it around. Is the spot a disk because that's the shape of the sun, or because it's the shape of my watch? In other words, would a square watch make a square spot, or do we just have a circular image of the circular sun, which will be circular no matter what? □



Problem 18.

15 Apply the equation $M = d_i/d_o$ to the case of a flat mirror.
▷ Solution, p. 1050 □

16 Use the method described in the text to derive the equation relating object distance to image distance for the case of a virtual image produced by a converging mirror. ▷ Solution, p. 1051 □

17 Find the focal length of the mirror in problem 6. ✓ □

18 Rank the focal lengths of the mirrors in the figure, from shortest to longest. Explain. □

19 (a) A converging mirror with a focal length of 20 cm is used to create an image, using an object at a distance of 10 cm. Is the image real, or is it virtual? (b) How about $f = 20$ cm and $d_o = 30$ cm? (c) What if it was a *diverging* mirror with $f = 20$ cm and $d_o = 10$ cm? (d) A diverging mirror with $f = 20$ cm and $d_o = 30$ cm?
▷ Solution, p. 1051 □

20 (a) Make up a numerical example of a virtual image formed by a converging mirror with a certain focal length, and determine the magnification. (You will need the result of problem 16.) Make sure to choose values of d_o and f that would actually produce a virtual image, not a real one. Now change the location of the object *a little bit* and redetermine the magnification, showing that it changes. At my local department store, the cosmetics department sells hand mirrors advertised as giving a magnification of 5 times. How would you interpret this?

(b) Suppose a Newtonian telescope is being used for astronomical observing. Assume for simplicity that no eyepiece is used, and assume a value for the focal length of the mirror that would be reasonable for an amateur instrument that is to fit in a closet. Is the angular magnification different for objects at different distances? For example, you could consider two planets, one of which is twice as far as the other.
▷ Solution, p. 1051 □

21 (a) Find a case where the magnification of a curved mirror is infinite. Is the *angular* magnification infinite from any realistic viewing position? (b) Explain why an arbitrarily large magnification can't be achieved by having a sufficiently small value of d_o .

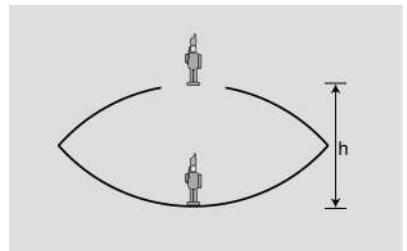
▷ Solution, p. 1051 □

22 A concave surface that reflects sound waves can act just like a converging mirror. Suppose that, standing near such a surface, you are able to find a point where you can place your head so that your own whispers are focused back on your head, so that they sound loud to you. Given your distance to the surface, what is the surface's focal length? ✓ ■

23 The figure shows a device for constructing a realistic optical illusion. Two mirrors of equal focal length are put against each other with their silvered surfaces facing inward. A small object placed in the bottom of the cavity will have its image projected in the air above. The way it works is that the top mirror produces a virtual image, and the bottom mirror then creates a real image of the virtual image. (a) Show that if the image is to be positioned as shown, at the mouth of the cavity, then the focal length of the mirrors is related to the dimension h via the equation

$$\frac{1}{f} = \frac{1}{h} + \frac{1}{h + \left(\frac{1}{h} - \frac{1}{f}\right)^{-1}}.$$

(b) Restate the equation in terms of a single variable $x = h/f$, and show that there are two solutions for x . Which solution is physically consistent with the assumptions of the calculation? ■



Problem 23.

24 (a) A converging mirror is being used to create a virtual image. What is the range of possible magnifications? (b) Do the same for the other types of images that can be formed by curved mirrors (both converging and diverging). ■

25 A diverging mirror of focal length f is fixed, and faces down. An object is dropped from the surface of the mirror, and falls away from it with acceleration g . The goal of the problem is to find the maximum velocity of the image.

- (a) Describe the motion of the image verbally, and explain why we should expect there to be a maximum velocity.
- (b) Use arguments based on units to determine the form of the solution, up to an unknown unitless multiplicative constant.
- (c) Complete the solution by determining the unitless constant.

✓ ■

26 Diamond has an index of refraction of 2.42, and part of the reason diamonds sparkle is that this encourages a light ray to undergo many total internal reflections before it emerges. (a) Calculate the critical angle at which total internal reflection occurs in diamond. (b) Explain the interpretation of your result: Is it measured from the normal, or from the surface? Is it a minimum angle for total internal reflection, or is it a maximum? How would the critical angle have been different for a substance such as glass or plastic, with a lower index of refraction? ✓ ■

27 Suppose a converging lens is constructed of a type of plastic whose index of refraction is less than that of water. How will the lens's behavior be different if it is placed underwater?

► Solution, p. 1052 ■

28 There are two main types of telescopes, refracting (using a lens) and reflecting (using a mirror as in figure i on p. 786). (Some telescopes use a mixture of the two types of elements: the light first encounters a large curved mirror, and then goes through an eyepiece that is a lens. To keep things simple, assume no eyepiece is used.) What implications would the color-dependence of focal length have for the relative merits of the two types of telescopes? Describe the case where an image is formed of a white star. You may find it helpful to draw a ray diagram. ■

29 Based on Snell's law, explain why rays of light passing through the edges of a converging lens are bent more than rays passing through parts closer to the center. It might seem like it should be the other way around, since the rays at the edge pass through less glass — shouldn't they be affected less? In your answer:

- Include a ray diagram showing a huge, full-page, close-up view of the relevant part of the lens.
- Make use of the fact that the front and back surfaces aren't always parallel; a lens in which the front and back surfaces *are* always parallel doesn't focus light at all, so if your explanation doesn't make use of this fact, your argument must be incorrect.
- Make sure your argument still works even if the rays don't come in parallel to the axis or from a point on the axis.

► Solution, p. 1052 ■

30 When you take pictures with a camera, the distance between the lens and the film or chip has to be adjusted, depending on the distance at which you want to focus. This is done by moving the lens. If you want to change your focus so that you can take a picture of something farther away, which way do you have to move the lens? Explain using ray diagrams. [Based on a problem by Eric Mazur.] ■

31 When swimming underwater, why is your vision made much clearer by wearing goggles with flat pieces of glass that trap air behind them? [Hint: You can simplify your reasoning by considering the special case where you are looking at an object far away, and along the optic axis of the eye.] \triangleright Solution, p. 1053 ■

32 An object is more than one focal length from a converging lens. (a) Draw a ray diagram. (b) Using reasoning like that developed in section 12.3, determine the positive and negative signs in the equation $1/f = \pm 1/d_i \pm 1/d_o$. (c) The images of the rose in section 4.2 were made using a lens with a focal length of 23 cm. If the lens is placed 80 cm from the rose, locate the image. \checkmark ■

33 The figure shows four lenses. Lens 1 has two spherical surfaces. Lens 2 is the same as lens 1 but turned around. Lens 3 is made by cutting through lens 1 and turning the bottom around. Lens 4 is made by cutting a central circle out of lens 1 and recessing it.

(a) A parallel beam of light enters lens 1 from the left, parallel to its axis. Reasoning based on Snell's law, will the beam emerging from the lens be bent inward, or outward, or will it remain parallel to the axis? Explain your reasoning. As part of your answer, make a huge drawing of one small part of the lens, and apply Snell's law at both interfaces. Recall that rays are bent more if they come to the interface at a larger angle with respect to the normal.

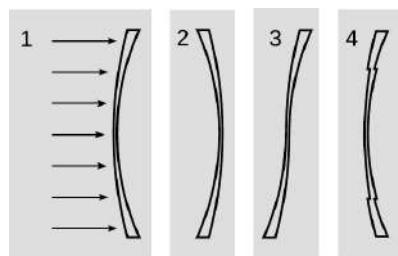
(b) What will happen with lenses 2, 3, and 4? Explain. Drawings are not necessary. \triangleright Solution, p. 1053 ■

34 The drawing shows the anatomy of the human eye, at twice life size. Find the radius of curvature of the outer surface of the cornea by measurements on the figure, and then derive the focal length of the air-cornea interface, where almost all the focusing of light occurs. You will need to use physical reasoning to modify the lensmaker's equation for the case where there is only a single refracting surface. Assume that the index of refraction of the cornea is essentially that of water.

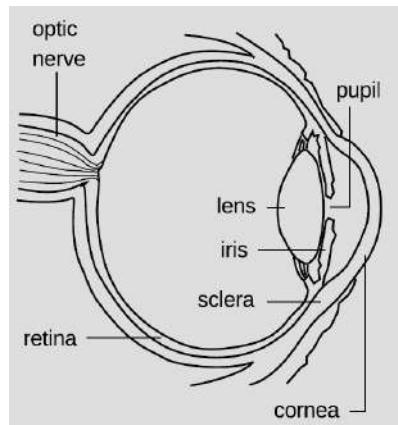
\checkmark ■

35 An object is less than one focal length from a converging lens. (a) Draw a ray diagram. (b) Using reasoning like that developed in section 12.3, determine the positive and negative signs in the equation $1/f = \pm 1/d_i \pm 1/d_o$. (c) The images of the rose in section 4.2 were made using a lens with a focal length of 23 cm. If the lens is placed 10 cm from the rose, locate the image. \checkmark ■

36 Nearsighted people wear glasses whose lenses are diverging. (a) Draw a ray diagram. For simplicity pretend that there is no eye behind the glasses. (b) Using reasoning like that developed in section 12.3, determine the positive and negative signs in the equation $1/f = \pm 1/d_i \pm 1/d_o$. (c) If the focal length of the lens is



Problem 33.



Problem 34.

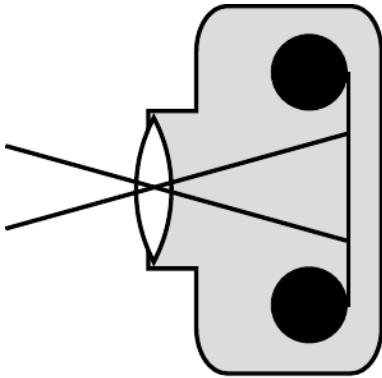
50.0 cm, and the person is looking at an object at a distance of 80.0 cm, locate the image. ✓ ■

37 (a) Light is being reflected diffusely from an object 1.000 m underwater. The light that comes up to the surface is refracted at the water-air interface. If the refracted rays all appear to come from the same point, then there will be a virtual image of the object in the water, above the object's actual position, which will be visible to an observer above the water. Consider three rays, A, B and C, whose angles in the water with respect to the normal are $\theta_i = 0.000^\circ$, 1.000° and 20.000° respectively. Find the depth of the point at which the refracted parts of A and B appear to have intersected, and do the same for A and C. Show that the intersections are at nearly the same depth, but not quite. [Check: The difference in depth should be about 4 cm.]

(b) Since all the refracted rays do not quite appear to have come from the same point, this is technically not a virtual image. In practical terms, what effect would this have on what you see?

(c) In the case where the angles are all small, use algebra and trig to show that the refracted rays do appear to come from the same point, and find an equation for the depth of the virtual image. Do not put in any numerical values for the angles or for the indices of refraction — just keep them as symbols. You will need the approximation $\sin \theta \approx \tan \theta \approx \theta$, which is valid for small angles measured in radians. ■

38 Prove that the principle of least time leads to Snell's law. ■



Problem 39.

39 Two standard focal lengths for camera lenses are 50 mm (standard) and 28 mm (wide-angle). To see how the focal lengths relate to the angular size of the field of view, it is helpful to visualize things as represented in the figure. Instead of showing many rays coming from the same point on the same object, as we normally do, the figure shows two rays from two different objects. Although the lens will intercept infinitely many rays from each of these points, we have shown only the ones that pass through the center of the lens, so that they suffer no angular deflection. (Any angular deflection at the front surface of the lens is canceled by an opposite deflection at the back, since the front and back surfaces are parallel at the lens's center.) What is special about these two rays is that they are aimed at the edges of one 35-mm-wide frame of film; that is, they show the limits of the field of view. Throughout this problem, we assume that d_o is much greater than d_i . (a) Compute the angular width of the camera's field of view when these two lenses are used. (b) Use small-angle approximations to find a simplified equation for the angular width of the field of view, θ , in terms of the focal length, f , and the width of the film, w . Your equation should not have any trig functions in it. Compare the results of this approximation with your answers from part a. (c) Suppose that we are holding

constant the aperture (amount of surface area of the lens being used to collect light). When switching from a 50-mm lens to a 28-mm lens, how many times longer or shorter must the exposure be in order to make a properly developed picture, i.e., one that is not under- or overexposed? [Based on a problem by Arnold Arons.]

▷ Solution, p. 1054 ■

40 A nearsighted person is one whose eyes focus light too strongly, and who is therefore unable to relax the lens inside her eye sufficiently to form an image on her retina of an object that is too far away.

(a) Draw a ray diagram showing what happens when the person tries, with uncorrected vision, to focus at infinity.

(b) What type of lenses do her glasses have? Explain.

(c) Draw a ray diagram showing what happens when she wears glasses. Locate both the image formed by the glasses and the final image.

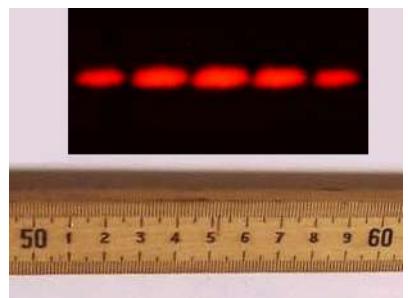
(d) Suppose she sometimes uses contact lenses instead of her glasses. Does the focal length of her contacts have to be less than, equal to, or greater than that of her glasses? Explain. ■

41 Fred's eyes are able to focus on things as close as 5.0 cm. Fred holds a magnifying glass with a focal length of 3.0 cm at a height of 2.0 cm above a flatworm. (a) Locate the image, and find the magnification. (b) Without the magnifying glass, from what distance would Fred want to view the flatworm to see its details as well as possible? With the magnifying glass? (c) Compute the angular magnification. ■

42 This problem has been deleted. 

43 It would be annoying if your eyeglasses produced a magnified or reduced image. Prove that when the eye is very close to a lens, and the lens produces a virtual image, the angular magnification is always approximately equal to 1 (regardless of whether the lens is diverging or converging). 

44 The figure shows a diffraction pattern made by a double slit, along with an image of a meter stick to show the scale. Sketch the diffraction pattern from the figure on your paper. Now consider the four variables in the equation $\lambda/d = \sin \theta/m$. Which of these are the same for all five fringes, and which are different for each fringe? Which variable would you naturally use in order to label which fringe was which? Label the fringes on your sketch using the values of that variable.

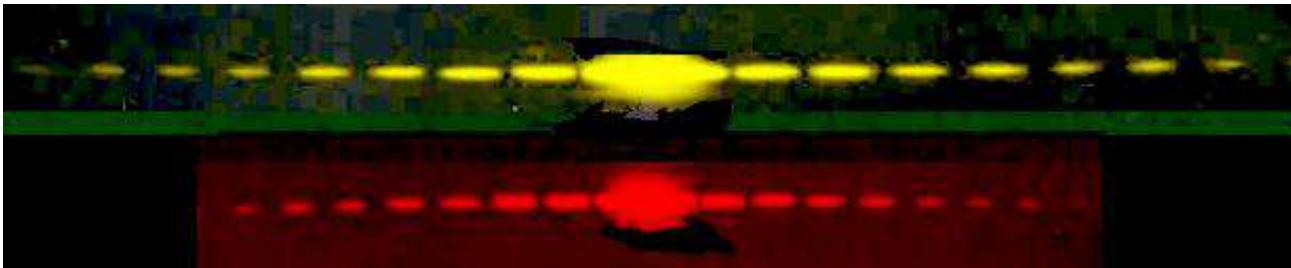


Problems 44 and 47.

45 Match gratings A-C with the diffraction patterns 1-3 that they produce. Explain.



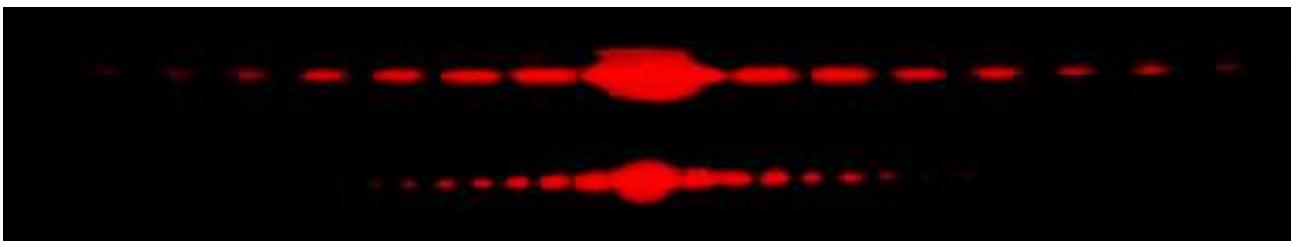
46 The figure below shows two diffraction patterns. The top one was made with yellow light, and the bottom one with red. Could the slits used to make the two patterns have been the same?



47 The figure on p. 839 shows a diffraction pattern made by a double slit, along with an image of a meter stick to show the scale. The slits were 146 cm away from the screen on which the diffraction pattern was projected. The spacing of the slits was 0.050 mm. What was the wavelength of the light? ✓ ■

48 Why would blue or violet light be the best for microscopy?
▷ Solution, p. 1054 ■

49 The figure below shows two diffraction patterns, both made with the same wavelength of red light. (a) What type of slits made the patterns? Is it a single slit, double slits, or something else? Explain. (b) Compare the dimensions of the slits used to make the top and bottom pattern. Give a numerical ratio, and state which way the ratio is, i.e., which slit pattern was the larger one. Explain.



▷ Solution, p. 1054 ■

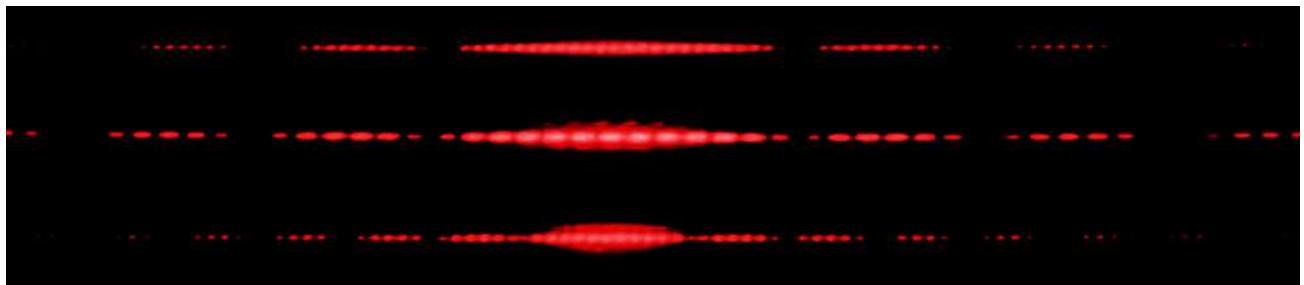
50 When white light passes through a diffraction grating, what is the smallest value of m for which the visible spectrum of order m overlaps the next one, of order $m + 1$? (The visible spectrum runs from about 400 nm to about 700 nm.) ■



Problem 51. This image of the Pleiades star cluster shows haloes around the stars due to the wave nature of light.

51 For star images such as the ones in figure y, estimate the angular width of the diffraction spot due to diffraction at the mouth of the telescope. Assume a telescope with a diameter of 10 meters (the largest currently in existence), and light with a wavelength in the middle of the visible range. Compare with the actual angular size of a star of diameter 10^9 m seen from a distance of 10^{17} m. What does this tell you? ▷ Solution, p. 1054 ■

52 The figure below shows three diffraction patterns. All were made under identical conditions, except that a different set of double slits was used for each one. The slits used to make the top pattern had a center-to-center separation $d = 0.50$ mm, and each slit was $w = 0.04$ mm wide. (a) Determine d and w for the slits used to make the pattern in the middle. (b) Do the same for the slits used to make the bottom pattern.



▷ Solution, p. 1055 ■

53 The beam of a laser passes through a diffraction grating, fans out, and illuminates a wall that is perpendicular to the original beam, lying at a distance of 2.0 m from the grating. The beam is produced by a helium-neon laser, and has a wavelength of 694.3 nm. The grating has 2000 lines per centimeter. (a) What is the distance on the wall between the central maximum and the maxima immediately to its right and left? (b) How much does your answer change when you use the small-angle approximations $\theta \approx \sin \theta \approx \tan \theta$? \checkmark

54 Ultrasound, i.e., sound waves with frequencies too high to be audible, can be used for imaging fetuses in the womb or for breaking up kidney stones so that they can be eliminated by the body. Consider the latter application. Lenses can be built to focus sound waves, but because the wavelength of the sound is not all that small compared to the diameter of the lens, the sound will not be concentrated exactly at the geometrical focal point. Instead, a diffraction pattern will be created with an intense central spot surrounded by fainter rings. About 85% of the power is concentrated within the central spot. The angle of the first minimum (surrounding the central spot) is given by $\sin \theta = \lambda/b$, where b is the diameter of the lens. This is similar to the corresponding equation for a single slit, but with a factor of 1.22 in front which arises from the circular shape of the aperture. Let the distance from the lens to the patient's kidney stone be $L = 20$ cm. You will want $f > 20$ kHz, so that the sound is inaudible. Find values of b and f that would result in a usable design, where the central spot is small enough to lie within a kidney stone 1 cm in diameter. \blacksquare

55 Under what circumstances could one get a mathematically undefined result by solving the double-slit diffraction equation for θ ? Give a physical interpretation of what would actually be observed.

▷ Solution, p. 1055 \blacksquare

56 When ultrasound is used for medical imaging, the frequency may be as high as 5-20 MHz. Another medical application of ultrasound is for therapeutic heating of tissues inside the body; here, the frequency is typically 1-3 MHz. What fundamental physical reasons could you suggest for the use of higher frequencies for imaging? \blacksquare

57 Suppose we have a polygonal room whose walls are mirrors, and there a pointlike light source in the room. In most such examples, every point in the room ends up being illuminated by the light source after some finite number of reflections. A difficult mathematical question, first posed in the middle of the last century, is whether it is ever possible to have an example in which the whole room is *not* illuminated. (Rays are assumed to be absorbed if they strike exactly at a vertex of the polygon, or if they pass exactly through the plane of a mirror.)

The problem was finally solved in 1995 by G.W. Tokarsky, who found an example of a room that was not illuminable from a certain point. Figure 57 shows a slightly simpler example found two years later by D. Castro. If a light source is placed at either of the locations shown with dots, the other dot remains unilluminated, although every other point is lit up. It is not straightforward to prove rigorously that Castro's solution has this property. However, the plausibility of the solution can be demonstrated as follows.

Suppose the light source is placed at the right-hand dot. Locate all the images formed by single reflections. Note that they form a regular pattern. Convince yourself that none of these images illuminates the left-hand dot. Because of the regular pattern, it becomes plausible that even if we form images of images, images of images of images, etc., none of them will ever illuminate the other dot.

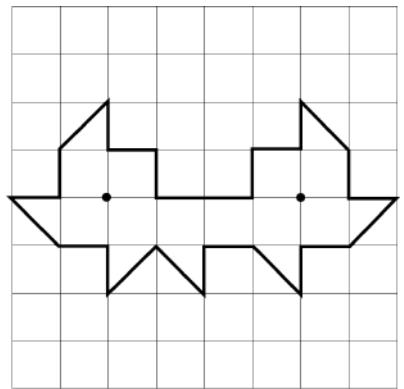
There are various other versions of the problem, some of which remain unsolved. The book by Klee and Wagon gives a good introduction to the topic, although it predates Tokarsky and Castro's work.

References:

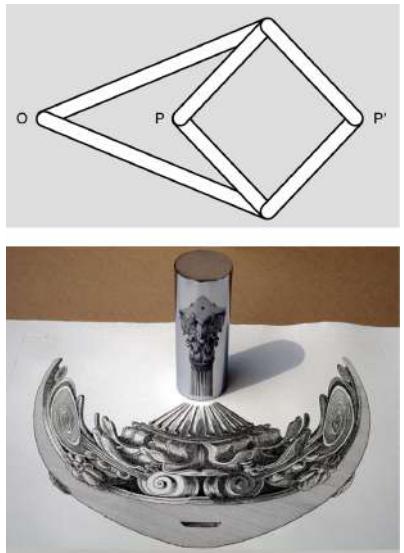
- G.W. Tokarsky, "Polygonal Rooms Not Illuminable from Every Point." Amer. Math. Monthly 102, 867-879, 1995.
 D. Castro, "Corrections." Quantum 7, 42, Jan. 1997.
 V. Klee and S. Wagon, *Old and New Unsolved Problems in Plane Geometry and Number Theory*. Mathematical Association of America, 1991. ■

58 A mechanical linkage is a device that changes one type of motion into another. The most familiar example occurs in a gasoline car's engine, where a connecting rod changes the linear motion of the piston into circular motion of the crankshaft. The top panel of the figure shows a mechanical linkage invented by Peaucellier in 1864, and independently by Lipkin around the same time. It consists of six rods joined by hinges, the four short ones forming a rhombus. Point O is fixed in space, but the apparatus is free to rotate about O. Motion at P is transformed into a different motion at P' (or vice versa).

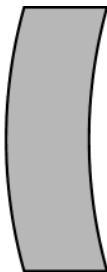
Geometrically, the linkage is a mechanical implementation of



Problem 57.



Problem 58.



Problem 59.

the ancient problem of inversion in a circle. Considering the case in which the rhombus is folded flat, let the k be the distance from O to the point where P and P' coincide. Form the circle of radius k with its center at O . As P and P' move in and out, points on the inside of the circle are always mapped to points on its outside, such that $rr' = k^2$. That is, the linkage is a type of analog computer that exactly solves the problem of finding the inverse of a number r . Inversion in a circle has many remarkable geometrical properties, discussed in H.S.M. Coxeter, *Introduction to Geometry*, Wiley, 1961. If a pen is inserted through a hole at P , and P' is traced over a geometrical figure, the Peaucellier linkage can be used to draw a kind of image of the figure.

A related problem is the construction of pictures, like the one in the bottom panel of the figure, called anamorphs. The drawing of the column on the paper is highly distorted, but when the reflecting cylinder is placed in the correct spot on top of the page, an undistorted image is produced inside the cylinder. (Wide-format movie technologies such as Cinemascope are based on similar principles.)

Show that the Peaucellier linkage does *not* convert correctly between an image and its anamorph, and design a modified version of the linkage that does. Some knowledge of analytic geometry will be helpful. ■

59 The figure shows a lens with surfaces that are curved, but whose thickness is constant along any horizontal line. Use the lens-maker's equation to prove that this "lens" is not really a lens at all. ▷ Solution, p. 1055 ■

60 Under ordinary conditions, gases have indices of refraction only a little greater than that of vacuum, i.e., $n = 1 + \epsilon$, where ϵ is some small number. Suppose that a ray crosses a boundary between a region of vacuum and a region in which the index of refraction is $1 + \epsilon$. Find the maximum angle by which such a ray can ever be deflected, in the limit of small ϵ . ▷ Hint, p. 1037 ✓ ■

61 A converging mirror has focal length f . An object is located at a distance $(1 + \epsilon)f$ from the mirror, where ϵ is small. Find the distance of the image from the mirror, simplifying your result as much as possible by using the assumption that ϵ is small.

▷ Answer, p. 1069 ■

62 The intensity of a beam of light is defined as the power per unit area incident on a perpendicular surface. Suppose that a beam of light in a medium with index of refraction n reaches the surface of the medium, with air on the outside. Its incident angle with respect to the normal is θ . (All angles are in radians.) Only a fraction f of the energy is transmitted, the rest being reflected. Because of this, we might expect that the transmitted ray would always be less intense than the incident one. But because the transmitted ray is refracted, it becomes narrower, causing an additional change in intensity by a factor $g > 1$. The product of these factors $I = fg$ can be greater than one. The purpose of this problem is to estimate the maximum amount of intensification.

We will use the small-angle approximation $\theta \ll 1$ freely, in order to make the math tractable. In our previous studies of waves, we have only studied the factor f in the one-dimensional case where $\theta = 0$. The generalization to $\theta \neq 0$ is rather complicated and depends on the polarization, but for unpolarized light, we can use Schlick's approximation,

$$f(\theta) = f(0)(1 - \cos \theta)^5,$$

where the value of f at $\theta = 0$ is found as in problem 17 on p. 395.

- (a) Using small-angle approximations, obtain an expression for g of the form $g \approx 1 + P\theta^2$, and find the constant P . \triangleright Answer, p. 1069
- (b) Find an expression for I that includes the two leading-order terms in θ . We will call this expression I_2 . Obtain a simple expression for the angle at which I_2 is maximized. As a check on your work, you should find that for $n = 1.3$, $\theta = 63^\circ$. (Trial-and-error maximization of I gives 60° .)
- (c) Find an expression for the maximum value of I_2 . You should find that for $n = 1.3$, the maximum intensification is 31%.



63 In an experiment to measure the unknown index of refraction n of a liquid, you send a laser beam from air into a tank filled with the liquid. Let ϕ be the angle of the beam relative to the normal while in the air, and let θ be the angle in the liquid. You can set ϕ to any value you like by aiming the laser from an appropriate direction, and you measure θ as a result. We wish to plan such an experiment so as to minimize the error dn in the result of the experiment, for a fixed error $d\theta$ in the measurement of the angle in the liquid. We assume that there is no significant contribution to the error from uncertainty in the index of refraction of air (which is very close to 1) or from the angle ϕ . Find dn in terms of $d\theta$, and determine the optimal conditions.

\triangleright Solution, p. 1055

64 Zahra likes to play practical jokes on the friends she goes hiking with. One night, by a blazing camp fire, she stealthily uses a lens of focal length f to gather light from the fire and make a hot spot on Becky's neck. (a) Using the method of section 12.3.2, p. 794, draw a ray diagram and set up the equation for the image location, inferring the correct plus and minus signs from the diagram. (b) Let A be the distance from the lens to the campfire, and B the distance from the lens to Becky's neck. Consider the following nine possibilities:

	B		
	$< f$	$= f$	$> f$
A	$< f$	<input type="checkbox"/>	<input type="checkbox"/>
	$= f$	<input type="checkbox"/>	<input type="checkbox"/>
	$> f$	<input type="checkbox"/>	<input type="checkbox"/>

By reasoning about your equation from part a, determine which of these are possible and which are not. \triangleright Solution, p. 1055 ■

Key to symbols:

■ easy ■ typical ■ challenging ■ difficult ■ very difficult
 ✓ An answer check is available at www.lightandmatter.com.

Exercises

Exercise 12A: Exploring Images With a Curved Mirror

Equipment:

- concave mirrors with deep curvature
- concave mirrors with gentle curvature
- convex mirrors

1. Obtain a curved mirror from your instructor. If it is silvered on both sides, make sure you're working with the concave side, which bends light rays inward. Look at your own face in the mirror. Now change the distance between your face and the mirror, and see what happens. Explore the full range of possible distances between your face and the mirror.

In these observations you've been changing two variables at once: the distance between the object (your face) and the mirror, and the distance from the mirror to your eye. In general, scientific experiments become easier to interpret if we practice isolation of variables, i.e., only change one variable while keeping all the others constant. In parts 2 and 3 you'll form an image of an object that's not your face, so that you can have independent control of the object distance and the point of view.

2. With the mirror held far away from you, observe the image of something behind you, over your shoulder. Now bring your eye closer and closer to the mirror. Can you see the image with your eye very close to the mirror? See if you can explain your observation by drawing a ray diagram.

—————> turn page

3. Now imagine the following new situation, but *don't actually do it yet*. Suppose you lay the mirror face-up on a piece of tissue paper, put your finger a few cm above the mirror, and look at the image of your finger. As in part 2, you can bring your eye closer and closer to the mirror.

Will you be able to see the image with your eye very close to the mirror? Draw a ray diagram to help you predict what you will observe.

Prediction: _____

Now test your prediction. If your prediction was incorrect, see if you can figure out what went wrong, or ask your instructor for help.

4. For parts 4 and 5, it's more convenient to use concave mirrors that are more gently curved; obtain one from your instructor. Lay the mirror on the tissue paper, and use it to create an image of the overhead lights on a piece of paper above it and a little off to the side. What do you have to do in order to make the image clear? Can you explain this observation using a ray diagram?

—————> turn page

5. Now imagine the following experiment, but *don't do it yet*. What will happen to the image on the paper if you cover half of the mirror with your hand?

Prediction: _____

Test your prediction. If your prediction was incorrect, can you explain what happened?

6. Now imagine forming an image with a convex mirror (one that bulges outward), and that therefore bends light rays away from the central axis (i.e., is diverging). Draw a typical ray diagram.

Is the image real, or virtual? Will there be more than one type of image?

Prediction: _____

Test your prediction.

Exercise 12B: Object and Image Distances

Equipment:

- optical benches
- converging mirrors
- illuminated objects

1. Set up the optical bench with the mirror at zero on the centimeter scale. Set up the illuminated object on the bench as well.
2. Each group will locate the image for their own value of the object distance, by finding where a piece of paper has to be placed in order to see the image on it. (The instructor will do one point as well.) Note that you will have to tilt the mirror a little so that the paper on which you project the image doesn't block the light from the illuminated object.

Is the image real or virtual? How do you know? Is it inverted, or uninverted?

Draw a ray diagram.

3. Measure the image distance and write your result in the table on the board. Do the same for the magnification.

4. What do you notice about the trend of the data on the board? Draw a second ray diagram with a different object distance, and show why this makes sense. Some tips for doing this correctly: (1) For simplicity, use the point on the object that is on the mirror's axis. (2) You need to trace two rays to locate the image. To save work, don't just do two rays at random angles. You can either use the on-axis ray as one ray, or do two rays that come off at the same angle, one above and one below the axis. (3) Where each ray hits the mirror, draw the normal line, and make sure the ray is at equal angles on both sides of the normal.

5. We will find the mirror's focal length from the instructor's data-point. Then, using this focal length, calculate a theoretical prediction of the image distance, and write it on the board next to the experimentally determined image distance.

Exercise 12C: How strong are your glasses?

This exercise was created by Dan MacIsaac.

Equipment:

eyeglasses

diverging lenses for students who don't wear glasses, or who use glasses with converging lenses

rulers and metersticks

scratch paper

marking pens

Most people who wear glasses have glasses whose lenses are outbending, which allows them to focus on objects far away. Such a lens cannot form a real image, so its focal length cannot be measured as easily as that of a converging lens. In this exercise you will determine the focal length of your own glasses by taking them off, holding them at a distance from your face, and looking through them at a set of parallel lines on a piece of paper. The lines will be reduced (the lens's magnification is less than one), and by adjusting the distance between the lens and the paper, you can make the magnification equal $1/2$ exactly, so that two spaces between lines as seen through the lens fit into one space as seen simultaneously to the side of the lens. This object distance can be used in order to find the focal length of the lens.

1. Use a marker to draw three evenly spaced parallel lines on the paper. (A spacing of a few cm works well.)
2. Does this technique really measure magnification or does it measure angular magnification? What can you do in your experiment in order to make these two quantities nearly the same, so the math is simpler?
3. Before taking any numerical data, use algebra to find the focal length of the lens in terms of d_o , the object distance that results in a magnification of $1/2$.
4. Measure the object distance that results in a magnification of $1/2$, and determine the focal length of your lens.

Exercise 12D: Double-Source Interference

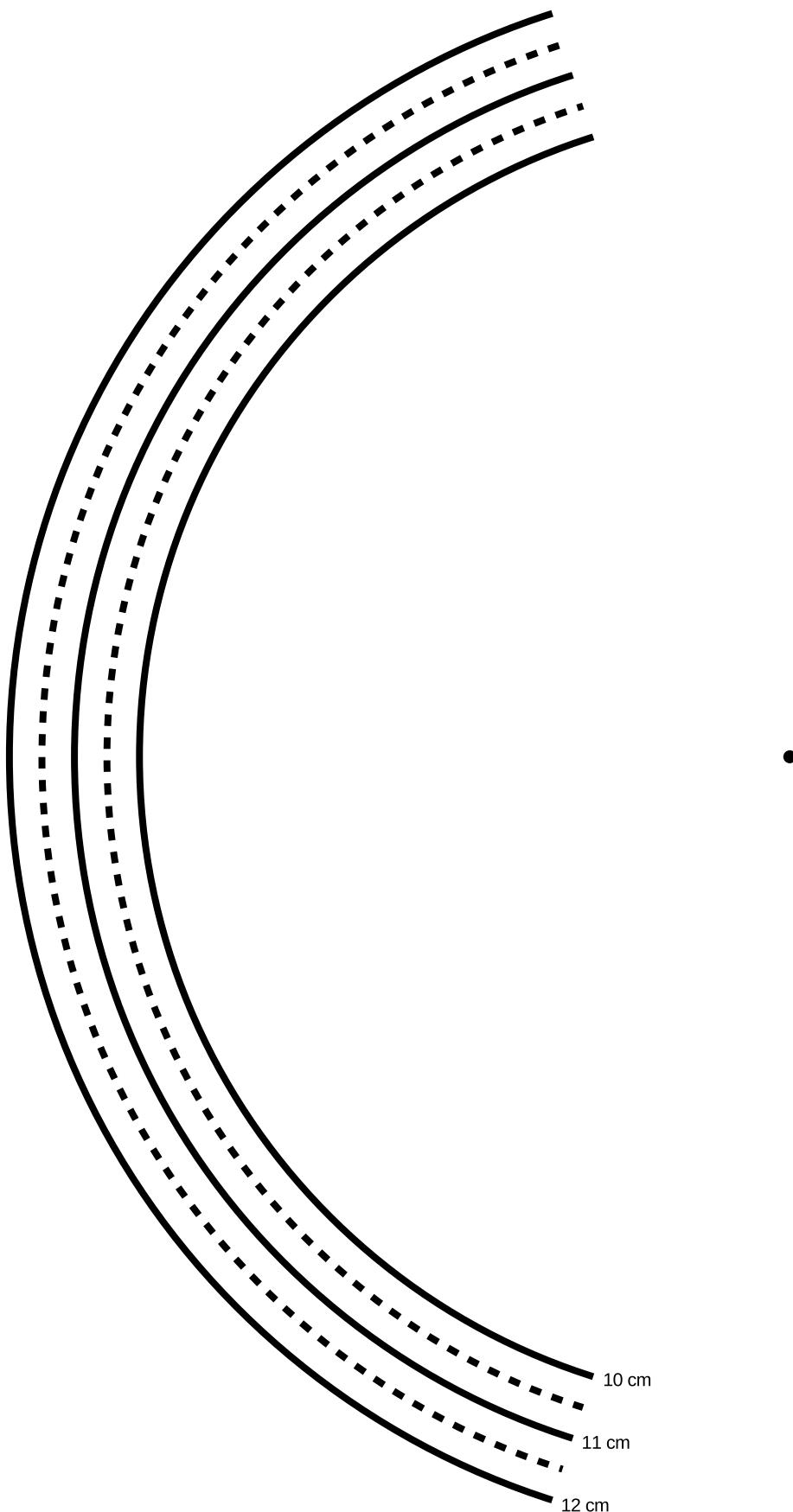
1. Two sources separated by a distance $d = 2$ cm make circular ripples with a wavelength of $\lambda = 1$ cm. On a piece of paper, make a life-size drawing of the two sources in the default setup, and locate the following points:

- A. The point that is 10 wavelengths from source #1 and 10 wavelengths from source #2.
- B. The point that is 10.5 wavelengths from #1 and 10.5 from #2.
- C. The point that is 11 wavelengths from #1 and 11 from #2.
- D. The point that is 10 wavelengths from #1 and 10.5 from #2.
- E. The point that is 11 wavelengths from #1 and 11.5 from #2.
- F. The point that is 10 wavelengths from #1 and 11 from #2.
- G. The point that is 11 wavelengths from #1 and 12 from #2.

You can do this either using a compass or by putting the next page under your paper and tracing. It is not necessary to trace all the arcs completely, and doing so is unnecessarily time-consuming; you can fairly easily estimate where these points would lie, and just trace arcs long enough to find the relevant intersections.

What do these points correspond to in the real wave pattern?

- 2. Make a fresh copy of your drawing, showing only point F and the two sources, which form a long, skinny triangle. Now suppose you were to change the setup by doubling d , while leaving λ the same. It's easiest to understand what's happening on the drawing if you move both sources outward, keeping the center fixed. Based on your drawing, what will happen to the position of point F when you double d ? How has the angle of point F changed?
- 3. What would happen if you doubled *both* λ and d compared to the standard setup?_____
- 4. Combining the ideas from parts 2 and 3, what do you think would happen to your angles if, starting from the standard setup, you doubled λ while leaving d the same?_____
- 5. Suppose λ was a millionth of a centimeter, while d was still as in the standard setup. What would happen to the angles? What does this tell you about observing diffraction of light?



Exercise 12E: Single-slit diffraction

Equipment:

rulers

computer with web browser

The following page is a diagram of a single slit and a screen onto which its diffraction pattern is projected. The class will make a numerical prediction of the intensity of the pattern at the different points on the screen. Each group will be responsible for calculating the intensity at one of the points. (Either 11 groups or six will work nicely – in the latter case, only points a, c, e, g, i, and k are used.) The idea is to break up the wavefront in the mouth of the slit into nine parts, each of which is assumed to radiate semicircular ripples as in Huygens' principle. The wavelength of the wave is 1 cm, and we assume for simplicity that each set of ripples has an amplitude of 1 unit when it reaches the screen.

1. For simplicity, let's imagine that we were only to use two sets of ripples rather than nine. You could measure the distance from each of the two points inside the slit to your point on the screen. Suppose the distances were both 25.0 cm. What would be the amplitude of the superimposed waves at this point on the screen?

Suppose one distance was 24.0 cm and the other was 25.0 cm. What would happen?

What if one was 24.0 cm and the other was 26.0 cm?

What if one was 24.5 cm and the other was 25.0 cm?

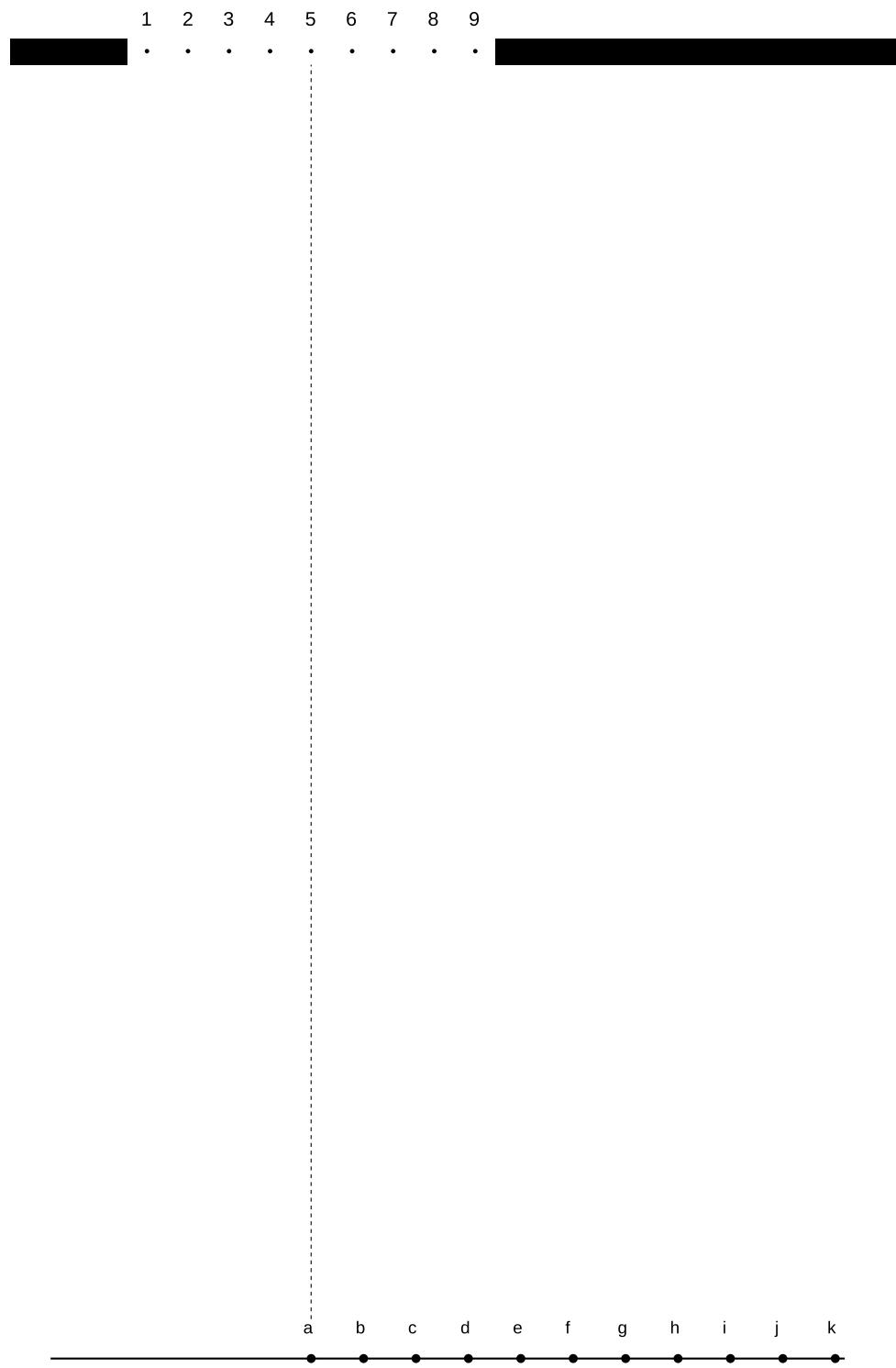
In general, what combinations of distances will lead to completely destructive and completely constructive interference?

Can you estimate the answer in the case where the distances are 24.7 and 25.0 cm?

2. Although it is possible to calculate mathematically the amplitude of the sine wave that results from superimposing two sine waves with an arbitrary phase difference between them, the algebra is rather laborious, and it becomes even more tedious when we have more than two waves to superimpose. Instead, one can simply use a computer spreadsheet or some other computer program to add up the sine waves numerically at a series of points covering one complete cycle. This is what we will actually do. You just need to enter the relevant data into the computer, then examine the results and pick off the amplitude from the resulting list of numbers. You can run the software through a web interface at <http://lightandmatter.com/cgi-bin/diffraction1.cgi>.

3. Measure all nine distances to your group's point on the screen, and write them on the board – that way everyone can see everyone else's data, and the class can try to make sense of why the results came out the way they did. Determine the amplitude of the combined wave, and write it on the board as well.

The class will discuss why the results came out the way they did.



Exercise 12F: Diffraction of Light

Equipment:

slit patterns, lasers, straight-filament bulbs

station 1

You have a mask with a bunch of different double slits cut out of it. The values of w and d are as follows:

pattern A	$w=0.04$ mm	$d=.250$ mm
pattern B	$w=0.04$ mm	$d=.500$ mm
pattern C	$w=0.08$ mm	$d=.250$ mm
pattern D	$w=0.08$ mm	$d=.500$ mm

Predict how the patterns will look different, and test your prediction. The easiest way to get the laser to point at different sets of slits is to stick folded up pieces of paper in one side or the other of the holders.

station 2

This is just like station 1, but with single slits:

pattern A	$w=0.02$ mm
pattern B	$w=0.04$ mm
pattern C	$w=0.08$ mm
pattern D	$w=0.16$ mm

Predict what will happen, and test your predictions. If you have time, check the actual numerical ratios of the w values against the ratios of the sizes of the diffraction patterns

station 3

This is like station 1, but the only difference among the sets of slits is how many slits there are:

pattern A	double slit
pattern B	3 slits
pattern C	4 slits
pattern D	5 slits

station 4

Hold the diffraction grating up to your eye, and look through it at the straight-filament light bulb. If you orient the grating correctly, you should be able to see the $m = 1$ and $m = -1$ diffraction patterns off the left and right. If you have it oriented the wrong way, they'll be above and below the bulb instead, which is inconvenient because the bulb's filament is vertical. Where is the $m = 0$ fringe? Can you see $m = 2$, etc.?

Station 5 has the same equipment as station 4. If you're assigned to station 5 first, you should actually do activity 4 first, because it's easier.

station 5

Use the transformer to increase and decrease the voltage across the bulb. This allows you to control the filament's temperature. Sketch graphs of intensity as a function of wavelength for various temperatures. The inability of the wave model of light to explain the mathematical shapes of these curves was historically one of the reasons for creating a new model, in which light is both a particle and a wave.

Chapter 13

Quantum Physics

13.1 Rules of randomness

Given for one instant an intelligence which could comprehend all the forces by which nature is animated and the respective positions of the things which compose it...nothing would be uncertain, and the future as the past would be laid out before its eyes.

Pierre Simon de Laplace, 1776

The energy produced by the atom is a very poor kind of thing. Anyone who expects a source of power from the transformation of these atoms is talking moonshine.

Ernest Rutherford, 1933

The Quantum Mechanics is very imposing. But an inner voice tells me that it is still not the final truth. The theory yields much, but it hardly brings us nearer to the secret of the Old One. In any case, I am convinced that He does not play dice.

Albert Einstein

However radical Newton's clockwork universe seemed to his contemporaries, by the early twentieth century it had become a sort of smugly accepted dogma. Luckily for us, this deterministic picture of the universe breaks down at the atomic level. The clearest demonstration that the laws of physics contain elements of randomness is in the behavior of radioactive atoms. Pick two identical atoms of a radioactive isotope, say the naturally occurring uranium 238, and watch them carefully. They will decay at different times, even though there was no difference in their initial behavior.

We would be in big trouble if these atoms' behavior was as predictable as expected in the Newtonian world-view, because radioactivity is an important source of heat for our planet. In reality, each atom chooses a random moment at which to release its energy, resulting in a nice steady heating effect. The earth would be a much colder planet if only sunlight heated it and not radioactivity. Probably there would be no volcanoes, and the oceans would never have been liquid. The deep-sea geothermal vents in which life first evolved would never have existed. But there would be an even worse consequence if radioactivity was deterministic: after a few billion years of peace, all the uranium 238 atoms in our planet would presumably pick the same moment to decay. The huge amount of stored nuclear



a / In 1980, the continental U.S. got its first taste of active volcanism in recent memory with the eruption of Mount St. Helens.

energy, instead of being spread out over eons, would all be released at one instant, blowing our whole planet to Kingdom Come.¹

The new version of physics, incorporating certain kinds of randomness, is called quantum physics (for reasons that will become clear later). It represented such a dramatic break with the previous, deterministic tradition that everything that came before is considered “classical,” even the theory of relativity. This chapter is a basic introduction to quantum physics.

Discussion Question

A I said “Pick two identical atoms of a radioactive isotope.” Are two atoms really identical? If their electrons are orbiting the nucleus, can we distinguish each atom by the particular arrangement of its electrons at some instant in time?

13.1.1 Randomness isn’t random.

Einstein’s distaste for randomness, and his association of determinism with divinity, goes back to the Enlightenment conception of the universe as a gigantic piece of clockwork that only had to be set in motion initially by the Builder. Many of the founders of quantum mechanics were interested in possible links between physics and Eastern and Western religious and philosophical thought, but every educated person has a different concept of religion and philosophy. Bertrand Russell remarked, “Sir Arthur Eddington deduces religion from the fact that atoms do not obey the laws of mathematics. Sir James Jeans deduces it from the fact that they do.”

Russell’s witticism, which implies incorrectly that mathematics cannot describe randomness, remind us how important it is not to oversimplify this question of randomness. You should not simply surmise, “Well, it’s all random, anything can happen.” For one thing, certain things simply cannot happen, either in classical physics or quantum physics. The conservation laws of mass, energy, momentum, and angular momentum are still valid, so for instance processes that create energy out of nothing are not just unlikely according to quantum physics, they are impossible.

A useful analogy can be made with the role of randomness in evolution. Darwin was not the first biologist to suggest that species changed over long periods of time. His two new fundamental ideas were that (1) the changes arose through random genetic variation, and (2) changes that enhanced the organism’s ability to survive and reproduce would be preserved, while maladaptive changes would be

¹This is under the assumption that all the uranium atoms were created at the same time. In reality, we have only a general idea of the processes that might have created the heavy elements in the nebula from which our solar system condensed. Some portion of them may have come from nuclear reactions in supernova explosions in that particular nebula, but some may have come from previous supernova explosions throughout our galaxy, or from exotic events like collisions of white dwarf stars.

eliminated by natural selection. Doubters of evolution often consider only the first point, about the randomness of natural variation, but not the second point, about the systematic action of natural selection. They make statements such as, “the development of a complex organism like *Homo sapiens* via random chance would be like a whirlwind blowing through a junkyard and spontaneously assembling a jumbo jet out of the scrap metal.” The flaw in this type of reasoning is that it ignores the deterministic constraints on the results of random processes. For an atom to violate conservation of energy is no more likely than the conquest of the world by chimpanzees next year.

Discussion Question

A Economists often behave like wannabe physicists, probably because it seems prestigious to make numerical calculations instead of talking about human relationships and organizations like other social scientists. Their striving to make economics work like Newtonian physics extends to a parallel use of mechanical metaphors, as in the concept of a market’s supply and demand acting like a self-adjusting machine, and the idealization of people as economic automatons who consistently strive to maximize their own wealth. What evidence is there for randomness rather than mechanical determinism in economics?

13.1.2 Calculating randomness

You should also realize that even if something is random, we can still understand it, and we can still calculate probabilities numerically. In other words, physicists are good bookmakers. A good bookmaker can calculate the odds that a horse will win a race much more accurately than an inexperienced one, but nevertheless cannot predict what will happen in any particular race.

Statistical independence

As an illustration of a general technique for calculating odds, suppose you are playing a 25-cent slot machine. Each of the three wheels has one chance in ten of coming up with a cherry. If all three wheels come up cherries, you win \$100. Even though the results of any particular trial are random, you can make certain quantitative predictions. First, you can calculate that your odds of winning on any given trial are $1/10 \times 1/10 \times 1/10 = 1/1000 = 0.001$. Here, I am representing the probabilities as numbers from 0 to 1, which is clearer than statements like “The odds are 999 to 1,” and makes the calculations easier. A probability of 0 represents something impossible, and a probability of 1 represents something that will definitely happen.

Also, you can say that any given trial is equally likely to result in a win, and it doesn’t matter whether you have won or lost in prior games. Mathematically, we say that each trial is statistically independent, or that separate games are uncorrelated. Most gamblers are mistakenly convinced that, to the contrary, games of chance are

correlated. If they have been playing a slot machine all day, they are convinced that it is “getting ready to pay,” and they do not want anyone else playing the machine and “using up” the jackpot that they “have coming.” In other words, they are claiming that a series of trials at the slot machine is negatively correlated, that losing now makes you more likely to win later. Craps players claim that you should go to a table where the person rolling the dice is “hot,” because she is likely to keep on rolling good numbers. Craps players, then, believe that rolls of the dice are positively correlated, that winning now makes you more likely to win later.

My method of calculating the probability of winning on the slot machine was an example of the following important rule for calculations based on independent probabilities:

the law of independent probabilities

If the probability of one event happening is P_A , and the probability of a second statistically independent event happening is P_B , then the probability that they will both occur is the product of the probabilities, $P_A P_B$. If there are more than two events involved, you simply keep on multiplying.

This can be taken as the definition of statistical independence.

Note that this only applies to independent probabilities. For instance, if you have a nickel and a dime in your pocket, and you randomly pull one out, there is a probability of 0.5 that it will be the nickel. If you then replace the coin and again pull one out randomly, there is again a probability of 0.5 of coming up with the nickel, because the probabilities are independent. Thus, there is a probability of 0.25 that you will get the nickel both times.

Suppose instead that you do not replace the first coin before pulling out the second one. Then you are bound to pull out the other coin the second time, and there is no way you could pull the nickel out twice. In this situation, the two trials are not independent, because the result of the first trial has an effect on the second trial. The law of independent probabilities does not apply, and the probability of getting the nickel twice is zero, not 0.25.

Experiments have shown that in the case of radioactive decay, the probability that any nucleus will decay during a given time interval is unaffected by what is happening to the other nuclei, and is also unrelated to how long it has gone without decaying. The first observation makes sense, because nuclei are isolated from each other at the centers of their respective atoms, and therefore have no physical way of influencing each other. The second fact is also reasonable, since all atoms are identical. Suppose we wanted to believe that certain atoms were “extra tough,” as demonstrated by their history of going an unusually long time without decaying. Those atoms would have to be different in some physical way, but nobody

has ever succeeded in detecting differences among atoms. There is no way for an atom to be changed by the experiences it has in its lifetime.

Addition of probabilities

The law of independent probabilities tells us to use multiplication to calculate the probability that both A and B will happen, assuming the probabilities are independent. What about the probability of an “or” rather than an “and”? If two events A and B are mutually exclusive, then the probability of one or the other occurring is the sum $P_A + P_B$. For instance, a bowler might have a 30% chance of getting a strike (knocking down all ten pins) and a 20% chance of knocking down nine of them. The bowler’s chance of knocking down either nine pins or ten pins is therefore 50%.

It does not make sense to add probabilities of things that are not mutually exclusive, i.e., that could both happen. Say I have a 90% chance of eating lunch on any given day, and a 90% chance of eating dinner. The probability that I will eat either lunch or dinner is not 180%.

Normalization

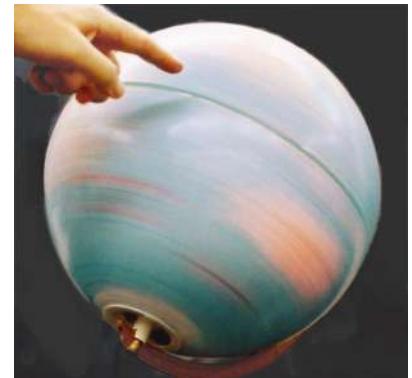
If I spin a globe and randomly pick a point on it, I have about a 70% chance of picking a point that’s in an ocean and a 30% chance of picking a point on land. The probability of picking either water or land is $70\% + 30\% = 100\%$. Water and land are mutually exclusive, and there are no other possibilities, so the probabilities had to add up to 100%. It works the same if there are more than two possibilities — if you can classify all possible outcomes into a list of mutually exclusive results, then all the probabilities have to add up to 1, or 100%. This property of probabilities is known as normalization.

Averages

Another way of dealing with randomness is to take averages. The casino knows that in the long run, the number of times you win will approximately equal the number of times you play multiplied by the probability of winning. In the slot-machine game described on page 859, where the probability of winning is 0.001, if you spend a week playing, and pay \$2500 to play 10,000 times, you are likely to win about 10 times ($10,000 \times 0.001 = 10$), and collect \$1000. On the average, the casino will make a profit of \$1500 from you. This is an example of the following rule.

Rule for Calculating Averages

If you conduct N identical, statistically independent trials, and the probability of success in each trial is P , then on the average, the total number of successful trials will be NP . If N is large enough, the relative error in this estimate will become



b / Normalization: the probability of picking land plus the probability of picking water adds up to 1.

small.

Foundations of probability

1. *Positivity*: Probabilities are positive.

2. *Normalization*: The total probability is 1.

3. *Additivity*: Mutually exclusive probabilities are additive.

4. *Independence*: Independent systems obey the definition of statistical independence, i.e., their probabilities multiply.

5. *The weak law of large numbers*: In the limit of a large number of trials, the frequency of a certain event converges to its probability.

Statements 1-3 are called the Kolmogorov axioms. In 3, for technical reasons, “additive” is usually taken to include infinite sums, such as $1 + 1/2 + 1/4 + \dots$, but not continuous sums such as integrals. Statement 4 is prediction about experiment. Statement 5 can be considered to be either a theorem to be proved from axioms such as 1-3, or an operational definition of probability, or a prediction about experiments.

The statement that the rule for calculating averages gets more and more accurate for larger and larger N (known popularly as the “law of averages”) often provides a correspondence principle that connects classical and quantum physics. For instance, the amount of power produced by a nuclear power plant is not random at any detectable level, because the number of atoms in the reactor is so large. In general, random behavior at the atomic level tends to average out when we consider large numbers of atoms, which is why physics seemed deterministic before physicists learned techniques for studying atoms individually.

We can achieve great precision with averages in quantum physics because we can use identical atoms to reproduce exactly the same situation many times. If we were betting on horses or dice, we would be much more limited in our precision. After a thousand races, the horse would be ready to retire. After a million rolls, the dice would be worn out.

When the number of trials is large, the accuracy of averages follows from the fact that the frequency of an event gets close to its probability. This is known as the law of large numbers.

The sidebar summarizes five basic facts that form the basis of probability theory.

self-check A

Which of the following things *must* be independent, which *could* be independent, and which definitely are *not* independent? (1) the probability of successfully making two free-throws in a row in basketball; (2) the probability that it will rain in London tomorrow and the probability that it will rain on the same day in a certain city in a distant galaxy; (3) your probability of dying today and of dying tomorrow. ▷ Answer, p. 1067

Discussion Questions

A Newtonian physics is an essentially perfect approximation for describing the motion of a pair of dice. If Newtonian physics is deterministic, why do we consider the result of rolling dice to be random?

B Why isn’t it valid to define randomness by saying that randomness is when all the outcomes are equally likely?

C The sequence of digits 1212121212121212 seems clearly nonrandom, and 41592653589793 seems random. The latter sequence, however, is the decimal form of pi, starting with the third digit. There is a story about the Indian mathematician Ramanujan, a self-taught prodigy, that a friend came to visit him in a cab, and remarked that the number of the cab, 1729, seemed relatively uninteresting. Ramanujan replied that on the contrary, it was very interesting because it was the smallest number that could be represented in two different ways as the sum of two cubes. The Argentine author Jorge Luis Borges wrote a short story called “The Library of Babel,” in which he imagined a library containing every book that could possibly be written using the letters of the alphabet. It would in-



c / Why are dice random?

clude a book containing only the repeated letter “a;” all the ancient Greek tragedies known today, all the lost Greek tragedies, and millions of Greek tragedies that were never actually written; your own life story, and various incorrect versions of your own life story; and countless anthologies containing a short story called “The Library of Babel.” Of course, if you picked a book from the shelves of the library, it would almost certainly look like a nonsensical sequence of letters and punctuation, but it’s always possible that the seemingly meaningless book would be a science-fiction screenplay written in the language of a Neanderthal tribe, or the lyrics to a set of incomparably beautiful love songs written in a language that never existed. In view of these examples, what does it really mean to say that something is random?

13.1.3 Probability distributions

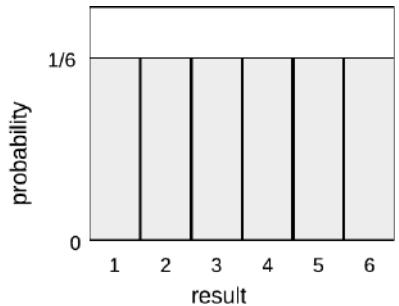
So far we’ve discussed random processes having only two possible outcomes: yes or no, win or lose, on or off. More generally, a random process could have a result that is a number. Some processes yield integers, as when you roll a die and get a result from one to six, but some are not restricted to whole numbers, for example the number of seconds that a uranium-238 atom will exist before undergoing radioactive decay.

Consider a throw of a die. If the die is “honest,” then we expect all six values to be equally likely. Since all six probabilities must add up to 1, then probability of any particular value coming up must be $1/6$. We can summarize this in a graph, d. Areas under the curve can be interpreted as total probabilities. For instance, the area under the curve from 1 to 3 is $1/6 + 1/6 + 1/6 = 1/2$, so the probability of getting a result from 1 to 3 is $1/2$. The function shown on the graph is called the probability distribution.

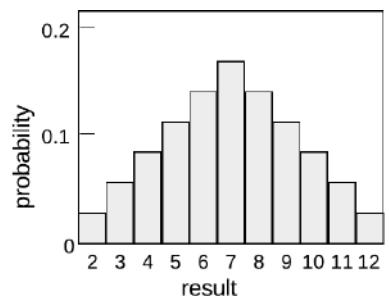
Figure e shows the probabilities of various results obtained by rolling two dice and adding them together, as in the game of craps. The probabilities are not all the same. There is a small probability of getting a two, for example, because there is only one way to do it, by rolling a one and then another one. The probability of rolling a seven is high because there are six different ways to do it: $1+6$, $2+5$, etc.

If the number of possible outcomes is large but finite, for example the number of hairs on a dog, the graph would start to look like a smooth curve rather than a ziggurat.

What about probability distributions for random numbers that are not integers? We can no longer make a graph with probability on the y axis, because the probability of getting a given exact number is typically zero. For instance, there is zero probability that a radioactive atom will last for *exactly* 3 seconds, since there are infinitely many possible results that are close to 3 but not exactly three: 2.999999999999996876876587658465436, for example. It doesn’t usually make sense, therefore, to talk about the probability



d / Probability distribution for the result of rolling a single die.

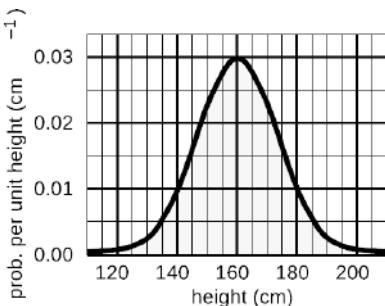


e / Rolling two dice and adding them up.

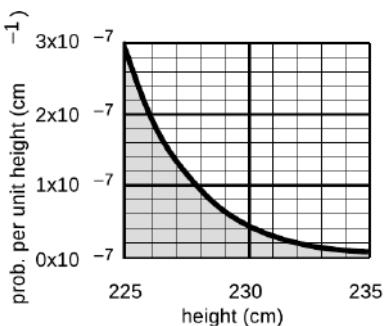
of a single numerical result, but it does make sense to talk about the probability of a certain range of results. For instance, the probability that an atom will last more than 3 and less than 4 seconds is a perfectly reasonable thing to discuss. We can still summarize the probability information on a graph, and we can still interpret areas under the curve as probabilities.

But the y axis can no longer be a unitless probability scale. In radioactive decay, for example, we want the x axis to have units of time, and we want areas under the curve to be unitless probabilities. The area of a single square on the graph paper is then

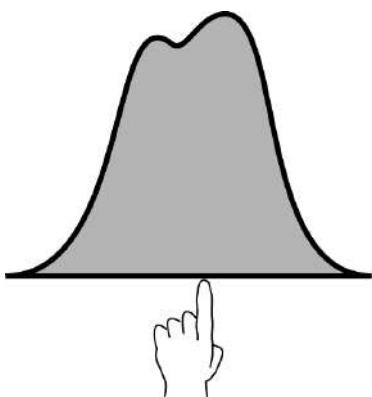
$$\begin{aligned} & \text{(unitless area of a square)} \\ & = (\text{width of square with time units}) \\ & \quad \times (\text{height of square}). \end{aligned}$$



f / A probability distribution for height of human adults (not real data).



g / Example 1.



h / The average of a probability distribution.

Compare the number of people with heights in the range of 130–135 cm to the number in the range 135–140.

▷ Answer, p. 1067

Looking for tall basketball players

example 1

▷ A certain country with a large population wants to find very tall people to be on its Olympic basketball team and strike a blow against western imperialism. Out of a pool of 10^8 people who are the right age and gender, how many are they likely to find who are over 225 cm (7 feet 4 inches) in height? Figure g gives a close-up of the “tail” of the distribution shown previously in figure f.

▷ The shaded area under the curve represents the probability that a given person is tall enough. Each rectangle represents a probability of $0.2 \times 10^{-7} \text{ cm}^{-1} \times 1 \text{ cm} = 2 \times 10^{-8}$. There are about 35 rectangles covered by the shaded area, so the probability of having a height greater than 225 cm is 7×10^{-7} , or just under one in a million. Using the rule for calculating averages, the average, or expected number of people this tall is $(10^8) \times (7 \times 10^{-7}) = 70$.

Average and width of a probability distribution

If the next Martian you meet asks you, “How tall is an adult human?,” you will probably reply with a statement about the average human height, such as “Oh, about 5 feet 6 inches.” If you wanted to explain a little more, you could say, “But that’s only an average. Most people are somewhere between 5 feet and 6 feet tall.” Without

bothering to draw the relevant bell curve for your new extraterrestrial acquaintance, you've summarized the relevant information by giving an average and a typical range of variation.

The average of a probability distribution can be defined geometrically as the horizontal position at which it could be balanced if it was constructed out of cardboard. A convenient numerical measure of the amount of variation about the average, or amount of uncertainty, is the full width at half maximum, or FWHM, shown in figure i.

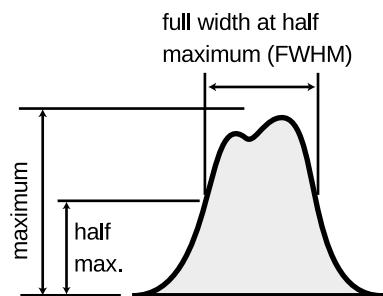
A great deal more could be said about this topic, and indeed an introductory statistics course could spend months on ways of defining the center and width of a distribution. Rather than force-feeding you on mathematical detail or techniques for calculating these things, it is perhaps more relevant to point out simply that there are various ways of defining them, and to inoculate you against the misuse of certain definitions.

The average is not the only possible way to say what is a typical value for a quantity that can vary randomly; another possible definition is the median, defined as the value that is exceeded with 50% probability. When discussing incomes of people living in a certain town, the average could be very misleading, since it can be affected massively if a single resident of the town is Bill Gates. Nor is the FWHM the only possible way of stating the amount of random variation; another possible way of measuring it is the standard deviation (defined as the square root of the average squared deviation from the average value).

13.1.4 Exponential decay and half-life

Half-life

Most people know that radioactivity “lasts a certain amount of time,” but that simple statement leaves out a lot. As an example, consider the following medical procedure used to diagnose thyroid function. A very small quantity of the isotope ^{131}I , produced in a nuclear reactor, is fed to or injected into the patient. The body’s biochemical systems treat this artificial, radioactive isotope exactly the same as ^{127}I , which is the only naturally occurring type. (Nutritionally, iodine is a necessary trace element. Iodine taken into the body is partly excreted, but the rest becomes concentrated in the thyroid gland. Iodized salt has had iodine added to it to prevent the nutritional deficiency known as goiters, in which the iodine-starved thyroid becomes swollen.) As the ^{131}I undergoes beta decay, it emits electrons, neutrinos, and gamma rays. The gamma rays can be measured by a detector passed over the patient’s body. As the radioactive iodine becomes concentrated in the thyroid, the amount of gamma radiation coming from the thyroid becomes greater, and that emitted by the rest of the body is reduced. The rate at which



i / The full width at half maximum (FWHM) of a probability distribution.

the iodine concentrates in the thyroid tells the doctor about the health of the thyroid.

If you ever undergo this procedure, someone will presumably explain a little about radioactivity to you, to allay your fears that you will turn into the Incredible Hulk, or that your next child will have an unusual number of limbs. Since iodine stays in your thyroid for a long time once it gets there, one thing you'll want to know is whether your thyroid is going to become radioactive forever. They may just tell you that the radioactivity "only lasts a certain amount of time," but we can now carry out a quantitative derivation of how the radioactivity really will die out.

Let $P_{surv}(t)$ be the probability that an iodine atom will survive without decaying for a period of at least t . It has been experimentally measured that half all ^{131}I atoms decay in 8 hours, so we have

$$P_{surv}(8 \text{ hr}) = 0.5.$$

Now using the law of independent probabilities, the probability of surviving for 16 hours equals the probability of surviving for the first 8 hours multiplied by the probability of surviving for the second 8 hours,

$$\begin{aligned} P_{surv}(16 \text{ hr}) &= 0.50 \times 0.50 \\ &= 0.25. \end{aligned}$$

Similarly we have

$$\begin{aligned} P_{surv}(24 \text{ hr}) &= 0.50 \times 0.5 \times 0.5 \\ &= 0.125. \end{aligned}$$

Generalizing from this pattern, the probability of surviving for any time t that is a multiple of 8 hours is

$$P_{surv}(t) = 0.5^{t/8 \text{ hr}}.$$

We now know how to find the probability of survival at intervals of 8 hours, but what about the points in time in between? What would be the probability of surviving for 4 hours? Well, using the law of independent probabilities again, we have

$$P_{surv}(8 \text{ hr}) = P_{surv}(4 \text{ hr}) \times P_{surv}(4 \text{ hr}),$$

which can be rearranged to give

$$\begin{aligned} P_{surv}(4 \text{ hr}) &= \sqrt{P_{surv}(8 \text{ hr})} \\ &= \sqrt{0.5} \\ &= 0.707. \end{aligned}$$

This is exactly what we would have found simply by plugging in $P_{surv}(t) = 0.5^{t/8 \text{ hr}}$ and ignoring the restriction to multiples of 8

hours. Since 8 hours is the amount of time required for half of the atoms to decay, it is known as the half-life, written $t_{1/2}$. The general rule is as follows:

Exponential Decay Equation

$$P_{\text{surv}}(t) = 0.5^{t/t_{1/2}}$$

Using the rule for calculating averages, we can also find the number of atoms, $N(t)$, remaining in a sample at time t :

$$N(t) = N(0) \times 0.5^{t/t_{1/2}}$$

Both of these equations have graphs that look like dying-out exponentials, as in the example below.

Radioactive contamination at Chernobyl *example 2*

- ▷ One of the most dangerous radioactive isotopes released by the Chernobyl disaster in 1986 was ${}^{90}\text{Sr}$, whose half-life is 28 years.
- (a) How long will it be before the contamination is reduced to one tenth of its original level? (b) If a total of 10^{27} atoms was released, about how long would it be before not a single atom was left?
- ▷ (a) We want to know the amount of time that a ${}^{90}\text{Sr}$ nucleus has a probability of 0.1 of surviving. Starting with the exponential decay formula,

$$P_{\text{surv}} = 0.5^{t/t_{1/2}},$$

we want to solve for t . Taking natural logarithms of both sides,

$$\ln P = \frac{t}{t_{1/2}} \ln 0.5,$$

so

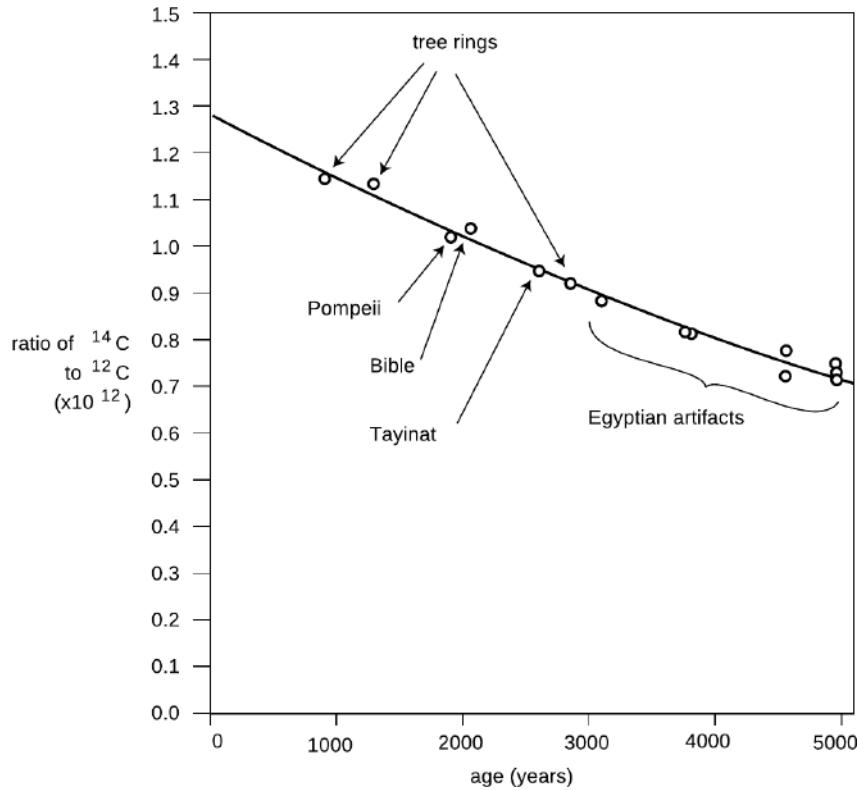
$$t = \frac{t_{1/2}}{\ln 0.5} \ln P$$

Plugging in $P = 0.1$ and $t_{1/2} = 28$ years, we get $t = 93$ years.

- (b) This is just like the first part, but $P = 10^{-27}$. The result is about 2500 years.

${}^{14}\text{C}$ Dating *example 3*

Almost all the carbon on Earth is ${}^{12}\text{C}$, but not quite. The isotope ${}^{14}\text{C}$, with a half-life of 5600 years, is produced by cosmic rays in the atmosphere. It decays naturally, but is replenished at such a rate that the fraction of ${}^{14}\text{C}$ in the atmosphere remains constant, at 1.3×10^{-12} . Living plants and animals take in both ${}^{12}\text{C}$ and ${}^{14}\text{C}$ from the atmosphere and incorporate both into their bodies. Once the living organism dies, it no longer takes in C atoms from the atmosphere, and the proportion of ${}^{14}\text{C}$ gradually falls off as it undergoes radioactive decay. This effect can be used to find the



j / Calibration of the ^{14}C dating method using tree rings and artifacts whose ages were known from other methods. Redrawn from Emilio Segrè, **Nuclei and Particles**, 1965.

age of dead organisms, or human artifacts made from plants or animals. Figure j on page 868 shows the exponential decay curve of ^{14}C in various objects. Similar methods, using longer-lived isotopes, provided the first firm proof that the earth was billions of years old, not a few thousand as some had claimed on religious grounds.

Rate of decay

If you want to find how many radioactive decays occur within a time interval lasting from time t to time $t + \Delta t$, the most straightforward approach is to calculate it like this:

$$\begin{aligned} & (\text{number of decays between } t \text{ and } t + \Delta t) \\ &= N(t) - N(t + \Delta t) \end{aligned}$$

Usually we're interested in the case where Δt is small compared to $t_{1/2}$, and in this limiting case the calculation starts to look exactly like the limit that goes into the definition of the derivative dN/dt . It is therefore more convenient to talk about the *rate* of decay – dN/dt rather than the *number* of decays in some finite time interval. Doing calculus on the function e^x is also easier than with 0.5^x , so we rewrite

the function $N(t)$ as

$$N = N(0)e^{-t/\tau},$$

where $\tau = t_{1/2}/\ln 2$ is shown in example 6 on p. 871 to be the average time of survival. The rate of decay is then

$$-\frac{dN}{dt} = \frac{N(0)}{\tau} e^{-t/\tau}.$$

Mathematically, differentiating an exponential just gives back another exponential. Physically, this is telling us that as N falls off exponentially, the rate of decay falls off at the same exponential rate, because a lower N means fewer atoms that remain available to decay.

self-check C

Check that both sides of the equation for the rate of decay have units of s^{-1} , i.e., decays per unit time. ▷ Answer, p. 1067

The hot potato

example 4

▷ A nuclear physicist with a demented sense of humor tosses you a cigar box, yelling “hot potato.” The label on the box says “contains 10^{20} atoms of ^{17}F , half-life of 66 s, produced today in our reactor at 1 p.m.” It takes you two seconds to read the label, after which you toss it behind some lead bricks and run away. The time is 1:40 p.m. Will you die?

▷ The time elapsed since the radioactive fluorine was produced in the reactor was 40 minutes, or 2400 s. The number of elapsed half-lives is therefore $t/t_{1/2} = 36$. The initial number of atoms was $N(0) = 10^{20}$. The number of decays per second is now about 10^7 s^{-1} , so it produced about 2×10^7 high-energy electrons while you held it in your hands. Although twenty million electrons sounds like a lot, it is not really enough to be dangerous.

By the way, none of the equations we’ve derived so far was the actual probability distribution for the time at which a particular radioactive atom will decay. That probability distribution would be found by substituting $N(0) = 1$ into the equation for the rate of decay.

Discussion Questions

A In the medical procedure involving ^{131}I , why is it the gamma rays that are detected, not the electrons or neutrinos that are also emitted?

B For 1 s, Fred holds in his hands 1 kg of radioactive stuff with a half-life of 1000 years. Ginger holds 1 kg of a different substance, with a half-life of 1 min, for the same amount of time. Did they place themselves in equal danger, or not?

C How would you interpret it if you calculated $N(t)$, and found it was less than one?

D Does the half-life depend on how much of the substance you have? Does the expected time until the sample decays completely depend on how much of the substance you have?

13.1.5 Applications of calculus

The area under the probability distribution is of course an integral. If we call the random number x and the probability distribution $D(x)$, then the probability that x lies in a certain range is given by

$$(\text{probability of } a \leq x \leq b) = \int_a^b D(x) dx.$$

What about averages? If x had a finite number of equally probable values, we would simply add them up and divide by how many we had. If they weren't equally likely, we'd make the weighted average $x_1 P_1 + x_2 P_2 + \dots$. But we need to generalize this to a variable x that can take on any of a continuum of values. The continuous version of a sum is an integral, so the average is

$$(\text{average value of } x) = \int x D(x) dx,$$

where the integral is over all possible values of x .

Probability distribution for radioactive decay *example 5*
 Here is a rigorous justification for the statement in subsection 13.1.4 that the probability distribution for radioactive decay is found by substituting $N(0) = 1$ into the equation for the rate of decay. We know that the probability distribution must be of the form

$$D(t) = k 0.5^{t/t_{1/2}},$$

where k is a constant that we need to determine. The atom is guaranteed to decay eventually, so normalization gives us

$$\begin{aligned} (\text{probability of } 0 \leq t < \infty) &= 1 \\ &= \int_0^\infty D(t) dt. \end{aligned}$$

The integral is most easily evaluated by converting the function into an exponential with e as the base

$$\begin{aligned} D(t) &= k \exp \left[\ln \left(0.5^{t/t_{1/2}} \right) \right] \\ &= k \exp \left[\frac{t}{t_{1/2}} \ln 0.5 \right] \\ &= k \exp \left(-\frac{\ln 2}{t_{1/2}} t \right), \end{aligned}$$

which gives an integral of the familiar form $\int e^{cx} dx = (1/c)e^{cx}$. We thus have

$$1 = -\frac{k t_{1/2}}{\ln 2} \exp \left(-\frac{\ln 2}{t_{1/2}} t \right) \Big|_0^\infty,$$

which gives the desired result:

$$k = \frac{\ln 2}{t_{1/2}}.$$

*Average lifetime**example 6*

You might think that the half-life would also be the average lifetime of an atom, since half the atoms' lives are shorter and half longer. But the half whose lives are longer include some that survive for many half-lives, and these rare long-lived atoms skew the average. We can calculate the average lifetime as follows:

$$(\text{average lifetime}) = \int_0^\infty t D(t) dt$$

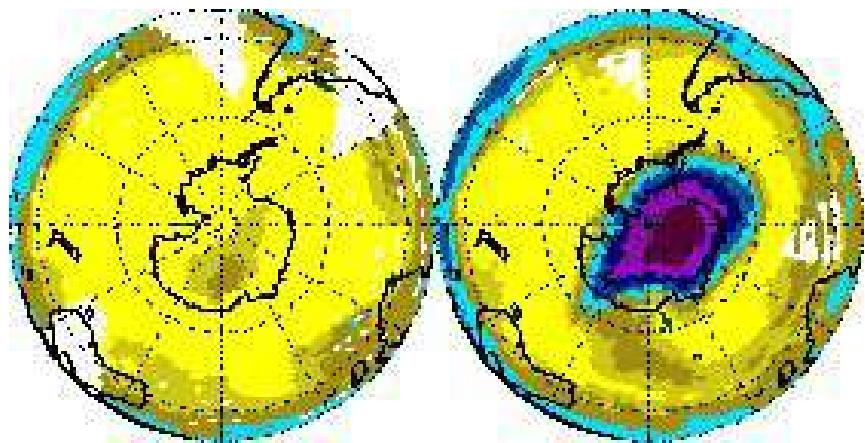
Using the convenient base-e form again, we have

$$(\text{average lifetime}) = \frac{\ln 2}{t_{1/2}} \int_0^\infty t \exp\left(-\frac{\ln 2}{t_{1/2}} t\right) dt.$$

This integral is of a form that can either be attacked with integration by parts or by looking it up in a table. The result is $\int xe^{cx} dx = \frac{x}{c} e^{cx} - \frac{1}{c^2} e^{cx}$, and the first term can be ignored for our purposes because it equals zero at both limits of integration. We end up with

$$\begin{aligned} (\text{average lifetime}) &= \frac{\ln 2}{t_{1/2}} \left(\frac{t_{1/2}}{\ln 2} \right)^2 \\ &= \frac{t_{1/2}}{\ln 2} \\ &= 1.443 t_{1/2}, \end{aligned}$$

which is, as expected, longer than one half-life.



k / In recent decades, a huge hole in the ozone layer has spread out from Antarctica. Left: November 1978. Right: November 1992

13.2 Light as a particle

The only thing that interferes with my learning is my education.
Albert Einstein

Radioactivity is random, but do the laws of physics exhibit randomness in other contexts besides radioactivity? Yes. Radioactive decay was just a good playpen to get us started with concepts of randomness, because all atoms of a given isotope are identical. By stocking the playpen with an unlimited supply of identical atom-toys, nature helped us to realize that their future behavior could be different regardless of their original identicality. We are now ready to leave the playpen, and see how randomness fits into the structure of physics at the most fundamental level.

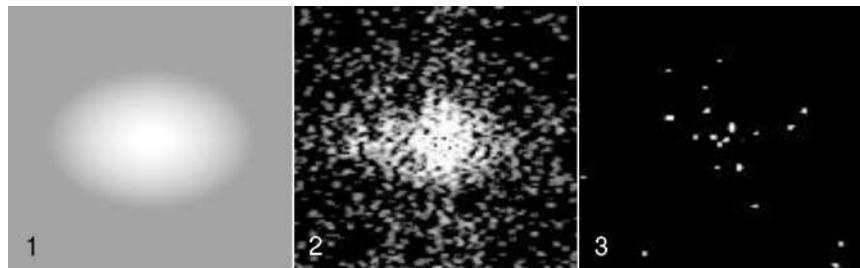
The laws of physics describe light and matter, and the quantum revolution rewrote both descriptions. Radioactivity was a good example of matter's behaving in a way that was inconsistent with classical physics, but if we want to get under the hood and understand how nonclassical things happen, it will be easier to focus on light rather than matter. A radioactive atom such as uranium-235 is after all an extremely complex system, consisting of 92 protons, 143 neutrons, and 92 electrons. Light, however, can be a simple sine wave.

However successful the classical wave theory of light had been — allowing the creation of radio and radar, for example — it still failed to describe many important phenomena. An example that is currently of great interest is the way the ozone layer protects us from the dangerous short-wavelength ultraviolet part of the sun's spectrum. In the classical description, light is a wave. When a wave

passes into and back out of a medium, its frequency is unchanged, and although its wavelength is altered while it is in the medium, it returns to its original value when the wave reemerges. Luckily for us, this is not at all what ultraviolet light does when it passes through the ozone layer, or the layer would offer no protection at all!

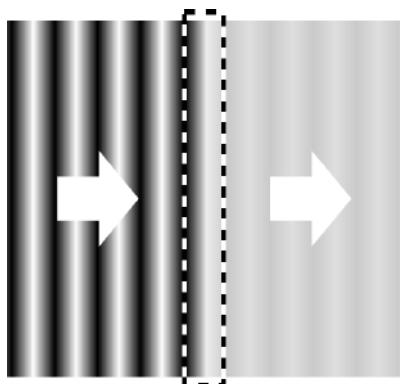
13.2.1 Evidence for light as a particle

For a long time, physicists tried to explain away the problems with the classical theory of light as arising from an imperfect understanding of atoms and the interaction of light with individual atoms and molecules. The ozone paradox, for example, could have been attributed to the incorrect assumption that one could think of the ozone layer as a smooth, continuous substance, when in reality it was made of individual ozone molecules. It wasn't until 1905 that Albert Einstein threw down the gauntlet, proposing that the problem had nothing to do with the details of light's interaction with atoms and everything to do with the fundamental nature of light itself.

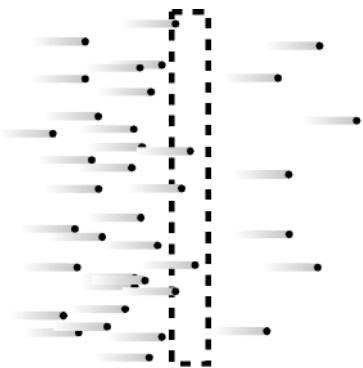


a / Digital camera images of dimmer and dimmer sources of light. The dots are records of individual photons.

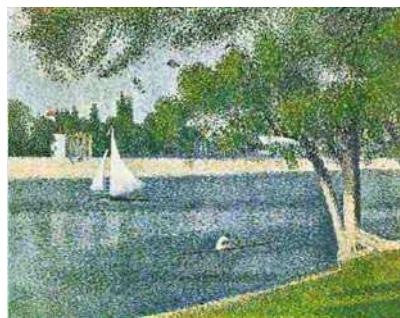
In those days the data were sketchy, the ideas vague, and the experiments difficult to interpret; it took a genius like Einstein to cut through the thicket of confusion and find a simple solution. Today, however, we can get right to the heart of the matter with a piece of ordinary consumer electronics, the digital camera. Instead of film, a digital camera has a computer chip with its surface divided up into a grid of light-sensitive squares, called "pixels." Compared to a grain of the silver compound used to make regular photographic film, a digital camera pixel is activated by an amount of light energy orders of magnitude smaller. We can learn something new about light by using a digital camera to detect smaller and smaller amounts of light, as shown in figure a. Figure a/1 is fake, but a/2 and a/3 are real digital-camera images made by Prof. Lyman Page of Princeton University as a classroom demonstration. Figure a/1 is what we would see if we used the digital camera to take a picture of a fairly



b / A wave is partially absorbed.



c / A stream of particles is partially absorbed.



d / Einstein and Seurat: twins separated at birth? *Seine Grande Jatte* by Georges Seurat (19th century).

dim source of light. In figures a/2 and a/3, the intensity of the light was drastically reduced by inserting semitransparent absorbers like the tinted plastic used in sunglasses. Going from a/1 to a/2 to a/3, more and more light energy is being thrown away by the absorbers.

The results are drastically different from what we would expect based on the wave theory of light. If light was a wave and nothing but a wave, b, then the absorbers would simply cut down the wave's amplitude across the whole wavefront. The digital camera's entire chip would be illuminated uniformly, and weakening the wave with an absorber would just mean that every pixel would take a long time to soak up enough energy to register a signal.

But figures a/2 and a/3 show that some pixels take strong hits while others pick up no energy at all. Instead of the wave picture, the image that is naturally evoked by the data is something more like a hail of bullets from a machine gun, c. Each “bullet” of light apparently carries only a tiny amount of energy, which is why detecting them individually requires a sensitive digital camera rather than an eye or a piece of film.

Although Einstein was interpreting different observations, this is the conclusion he reached in his 1905 paper: that the pure wave theory of light is an oversimplification, and that the energy of a beam of light comes in finite chunks rather than being spread smoothly throughout a region of space.

We now think of these chunks as particles of light, and call them “photons,” although Einstein avoided the word “particle,” and the word “photon” was invented later. Regardless of words, the trouble was that waves and particles seemed like inconsistent categories. The reaction to Einstein’s paper could be kindly described as vigorously skeptical. Even twenty years later, Einstein wrote, “There are therefore now two theories of light, both indispensable, and — as one must admit today despite twenty years of tremendous effort on the part of theoretical physicists — without any logical connection.” In the remainder of this section we will learn how the seeming paradox was eventually resolved.

Discussion Questions

A Suppose someone rebuts the digital camera data in figure a, claiming that the random pattern of dots occurs not because of anything fundamental about the nature of light but simply because the camera's pixels are not all exactly the same — some are just more sensitive than others. How could we test this interpretation?

B Discuss how the correspondence principle applies to the observations and concepts discussed in this section.

13.2.2 How much light is one photon?

The photoelectric effect

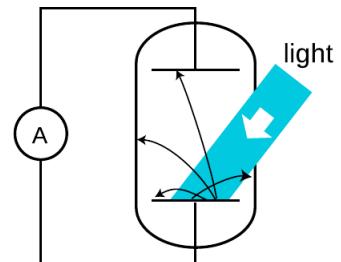
We have seen evidence that light energy comes in little chunks, so the next question to be asked is naturally how much energy is in one chunk. The most straightforward experimental avenue for addressing this question is a phenomenon known as the photoelectric effect. The photoelectric effect occurs when a photon strikes the surface of a solid object and knocks out an electron. It occurs continually all around you. It is happening right now at the surface of your skin and on the paper or computer screen from which you are reading these words. It does not ordinarily lead to any observable electrical effect, however, because on the average free electrons are wandering back in just as frequently as they are being ejected. (If an object did somehow lose a significant number of electrons, its growing net positive charge would begin attracting the electrons back more and more strongly.)

Figure e shows a practical method for detecting the photoelectric effect. Two very clean parallel metal plates (the electrodes of a capacitor) are sealed inside a vacuum tube, and only one plate is exposed to light. Because there is a good vacuum between the plates, any ejected electron that happens to be headed in the right direction will almost certainly reach the other capacitor plate without colliding with any air molecules.

The illuminated (bottom) plate is left with a net positive charge, and the unilluminated (top) plate acquires a negative charge from the electrons deposited on it. There is thus an electric field between the plates, and it is because of this field that the electrons' paths are curved, as shown in the diagram. However, since vacuum is a good insulator, any electrons that reach the top plate are prevented from responding to the electrical attraction by jumping back across the gap. Instead they are forced to make their way around the circuit, passing through an ammeter. The ammeter allows a measurement of the strength of the photoelectric effect.

An unexpected dependence on frequency

The photoelectric effect was discovered serendipitously by Heinrich Hertz in 1887, as he was experimenting with radio waves. He was not particularly interested in the phenomenon, but he did notice that the effect was produced strongly by ultraviolet light and more weakly by lower frequencies. Light whose frequency was lower than a certain critical value did not eject any electrons at all. (In fact this was all prior to Thomson's discovery of the electron, so Hertz would not have described the effect in terms of electrons — we are discussing everything with the benefit of hindsight.) This dependence on frequency didn't make any sense in terms of the classical wave theory of light. A light wave consists of electric and magnetic



e / Apparatus for observing the photoelectric effect. A beam of light strikes a capacitor plate inside a vacuum tube, and electrons are ejected (black arrows).

fields. The stronger the fields, i.e., the greater the wave's amplitude, the greater the forces that would be exerted on electrons that found themselves bathed in the light. It should have been amplitude (brightness) that was relevant, not frequency. The dependence on frequency not only proves that the wave model of light needs modifying, but with the proper interpretation it allows us to determine how much energy is in one photon, and it also leads to a connection between the wave and particle models that we need in order to reconcile them.



f / The hamster in her hamster ball is like an electron emerging from the metal (tiled kitchen floor) into the surrounding vacuum (wood floor). The wood floor is higher than the tiled floor, so as she rolls up the step, the hamster will lose a certain amount of kinetic energy, analogous to E_s . If her kinetic energy is too small, she won't even make it up the step.

To make any progress, we need to consider the physical process by which a photon would eject an electron from the metal electrode. A metal contains electrons that are free to move around. Ordinarily, in the interior of the metal, such an electron feels attractive forces from atoms in every direction around it. The forces cancel out. But if the electron happens to find itself at the surface of the metal, the attraction from the interior side is not balanced out by any attraction from outside. In popping out through the surface the electron therefore loses some amount of energy E_s , which depends on the type of metal used.

Suppose a photon strikes an electron, annihilating itself and giving up all its energy to the electron. (We now know that this is what always happens in the photoelectric effect, although it had not yet been established in 1905 whether or not the photon was completely annihilated.) The electron will (1) lose kinetic energy through collisions with other electrons as it plows through the metal on its way to the surface; (2) lose an amount of kinetic energy equal to E_s as it emerges through the surface; and (3) lose more energy on its way across the gap between the plates, due to the electric field between the plates. Even if the electron happens to be right at the surface of the metal when it absorbs the photon, and even if the electric field between the plates has not yet built up very much, E_s is the bare minimum amount of energy that it must receive from the photon if it is to contribute to a measurable current. The reason for using very clean electrodes is to minimize E_s and make it have a definite value characteristic of the metal surface, not a mixture of values due to the various types of dirt and crud that are present in tiny amounts on all surfaces in everyday life.

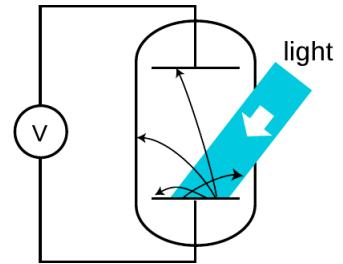
We can now interpret the frequency dependence of the photoelectric effect in a simple way: apparently the amount of energy possessed by a photon is related to its frequency. A low-frequency red or infrared photon has an energy less than E_s , so a beam of them will not produce any current. A high-frequency blue or violet photon, on the other hand, packs enough of a punch to allow an electron to make it to the other plate. At frequencies higher than the minimum, the photoelectric current continues to increase with the frequency of the light because of effects (1) and (3).

Numerical relationship between energy and frequency

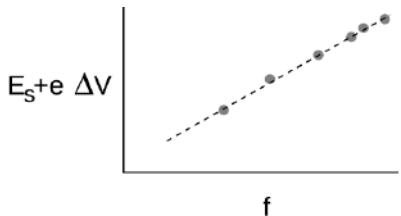
Figure g shows an experiment that is used sometimes in college laboratory courses to probe the relationship between the energy and frequency of a photon. The idea is simply to illuminate one plate of the vacuum tube with light of a single wavelength and monitor the voltage difference between the two plates as they charge up. Since the resistance of a voltmeter is very high (much higher than the resistance of an ammeter), we can assume to a good approximation that electrons reaching the top plate are stuck there permanently, so the voltage will keep on increasing for as long as electrons are making it across the vacuum tube.

At a moment when the voltage difference has reached a value ΔV , the minimum energy required by an electron to make it out of the bottom plate and across the gap to the other plate is $E_s + e\Delta V$. As ΔV increases, we eventually reach a point at which $E_s + e\Delta V$ equals the energy of one photon. No more electrons can cross the gap, and the reading on the voltmeter stops rising. The quantity $E_s + e\Delta V$ now tells us the energy of one photon. If we determine this energy for a variety of wavelengths, h , we find the following simple relationship between the energy of a photon and the frequency of the light:

$$E = hf,$$



g / A different way of studying the photoelectric effect.



h / The quantity $E_s + e\Delta V$ indicates the energy of one photon. It is found to be proportional to the frequency of the light.

where h is a constant with the value $6.63 \times 10^{-34} \text{ J} \cdot \text{s}$. Note how the equation brings the wave and particle models of light under the same roof: the left side is the energy of one *particle* of light, while the right side is the frequency of the same light, interpreted as a *wave*. The constant h is known as Planck's constant, for historical reasons explained in the footnote beginning on the preceding page.

self-check D

How would you extract h from the graph in figure h? What if you didn't even know E_s in advance, and could only graph $e\Delta V$ versus f ? ▶

Answer, p. 1067

Since the energy of a photon is hf , a beam of light can only have energies of hf , $2hf$, $3hf$, etc. Its energy is quantized — there is no such thing as a fraction of a photon. Quantum physics gets its name from the fact that it quantizes quantities like energy, momentum, and angular momentum that had previously been thought to be smooth, continuous and infinitely divisible.

Photons from a lightbulb

example 7

▷ Roughly how many photons are emitted by a 100 watt lightbulb in 1 second?

▷ People tend to remember wavelengths rather than frequencies for visible light. The bulb emits photons with a range of frequencies and wavelengths, but let's take 600 nm as a typical wavelength for purposes of estimation. The energy of a single photon is

$$\begin{aligned}E_{\text{photon}} &= hf \\&= hc/\lambda\end{aligned}$$

A power of 100 W means 100 joules per second, so the number of photons is

$$\begin{aligned}(100 \text{ J})/E_{\text{photon}} &= (100 \text{ J})/(hc/\lambda) \\&\approx 3 \times 10^{20}\end{aligned}$$

This hugeness of this number is consistent with the correspondence principle. The experiments that established the classical theory of optics weren't wrong. They were right, within their domain of applicability, in which the number of photons was so large as to be indistinguishable from a continuous beam.

Measuring the wave

example 8

When surfers are out on the water waiting for their chance to catch a wave, they're interested in both the height of the waves and when the waves are going to arrive. In other words, they observe both the amplitude and phase of the waves, and it doesn't matter to them that the water is granular at the molecular level. The correspondence principle requires that we be able to do the same thing for electromagnetic waves, since the classical theory of electricity and magnetism was all stated and verified experimentally in terms of the fields **E** and **B**, which are the amplitude of an electromagnetic wave. The phase is also necessary, since the induction effects predicted by Maxwell's equation would flip their signs depending on whether an oscillating field is on its way up or on its way back down.

This is a more demanding application of the correspondence principle than the one in example 7, since amplitudes and phases constitute more detailed information than the over-all intensity of a beam of light. Eyeball measurements can't detect this type of information, since the eye is much bigger than a wavelength, but for example an AM radio receiver can do it with radio waves, since the wavelength for a station at 1000 kHz is about 300 meters, which is much larger than the antenna. The correspondence principle demands that we be able to explain this in terms of the photon theory, and this requires not just that we have a large number

of photons emitted by the transmitter per second, as in example 7, but that even by the time they spread out and reach the receiving antenna, there should be many photons overlapping each other within a space of one cubic wavelength. Problem 47 on p. 952 verifies that the number is in fact extremely large.

Momentum of a photon

example 9

▷ According to the theory of relativity, the momentum of a beam of light is given by $p = E/c$. Apply this to find the momentum of a single photon in terms of its frequency, and in terms of its wavelength.

▷ Combining the equations $p = E/c$ and $E = hf$, we find

$$\begin{aligned} p &= E/c \\ &= \frac{h}{c}f. \end{aligned}$$

To reexpress this in terms of wavelength, we use $c = f\lambda$:

$$\begin{aligned} p &= \frac{h}{c} \cdot \frac{c}{\lambda} \\ &= \frac{h}{\lambda} \end{aligned}$$

The second form turns out to be simpler.

Discussion Questions

A The photoelectric effect only ever ejects a very tiny percentage of the electrons available near the surface of an object. How well does this agree with the wave model of light, and how well with the particle model? Consider the two different distance scales involved: the wavelength of the light, and the size of an atom, which is on the order of 10^{-10} or 10^{-9} m.

B What is the significance of the fact that Planck's constant is numerically very small? How would our everyday experience of light be different if it was not so small?

C How would the experiments described above be affected if a single electron was likely to get hit by more than one photon?

D Draw some representative trajectories of electrons for $\Delta V = 0$, ΔV less than the maximum value, and ΔV greater than the maximum value.

E Explain based on the photon theory of light why ultraviolet light would be more likely than visible or infrared light to cause cancer by damaging DNA molecules. How does this relate to discussion question C?

F Does $E = hf$ imply that a photon changes its energy when it passes from one transparent material into another substance with a different index of refraction?

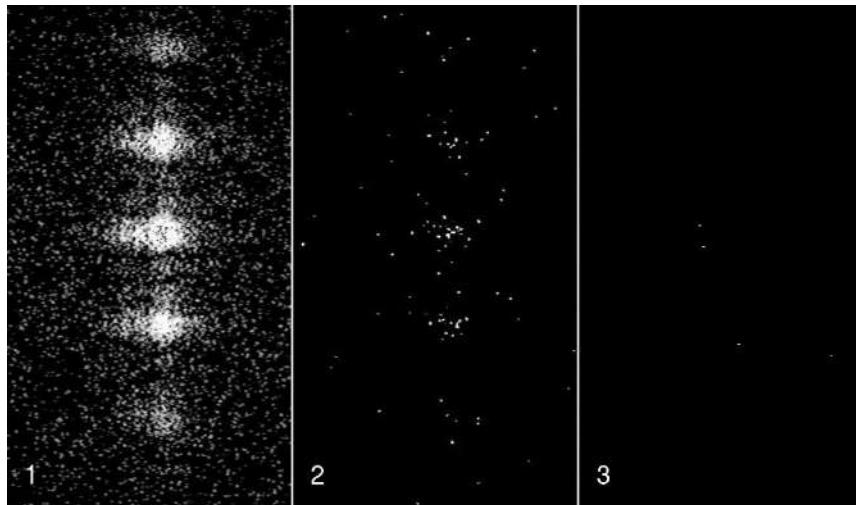
13.2.3 Wave-particle duality

How can light be both a particle and a wave? We are now ready to resolve this seeming contradiction. Often in science when something seems paradoxical, it's because we (1) don't define our terms carefully, or (2) don't test our ideas against any specific real-world situation. Let's define particles and waves as follows:

- Waves exhibit superposition, and specifically interference phenomena.
- Particles can only exist in whole numbers, not fractions.

As a real-world check on our philosophizing, there is one particular experiment that works perfectly. We set up a double-slit interference experiment that we know will produce a diffraction pattern if light is an honest-to-goodness wave, but we detect the light with a detector that is capable of sensing individual photons, e.g., a digital camera. To make it possible to pick out individual dots due to individual photons, we must use filters to cut down the intensity of the light to a very low level, just as in the photos by Prof. Page on p. 873. The whole thing is sealed inside a light-tight box. The results are shown in figure i. (In fact, the similar figures in on page 873 are simply cutouts from these figures.)

i / Wave interference patterns photographed by Prof. Lyman Page with a digital camera. Laser light with a single well-defined wavelength passed through a series of absorbers to cut down its intensity, then through a set of slits to produce interference, and finally into a digital camera chip. (A triple slit was actually used, but for conceptual simplicity we discuss the results in the main text as if it was a double slit.) In panel 2 the intensity has been reduced relative to 1, and even more so for panel 3.



Neither the pure wave theory nor the pure particle theory can explain the results. If light was only a particle and not a wave, there would be no interference effect. The result of the experiment would be like firing a hail of bullets through a double slit, j. Only two spots directly behind the slits would be hit.

If, on the other hand, light was only a wave and not a particle, we would get the same kind of diffraction pattern that would happen

with a water wave, k . There would be no discrete dots in the photo, only a diffraction pattern that shaded smoothly between light and dark.

Applying the definitions to this experiment, light must be both a particle and a wave. It is a wave because it exhibits interference effects. At the same time, the fact that the photographs contain discrete dots is a direct demonstration that light refuses to be split into units of less than a single photon. There can only be whole numbers of photons: four photons in figure i/3, for example.

A wrong interpretation: photons interfering with each other

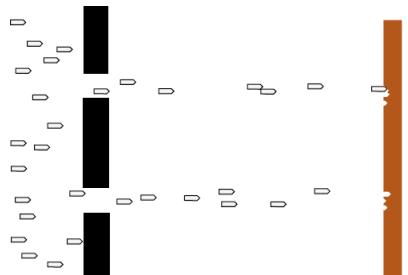
One possible interpretation of wave-particle duality that occurred to physicists early in the game was that perhaps the interference effects came from photons interacting with each other. By analogy, a water wave consists of moving water molecules, and interference of water waves results ultimately from all the mutual pushes and pulls of the molecules. This interpretation has been conclusively disproved by forming interference patterns with light so dim that no more than one photon is in flight at a time. In figure i/3, for example, the intensity of the light has been cut down so much by the absorbers that if it was in the open, the average separation between photons would be on the order of a kilometer! Although most light sources tend to emit photons in bunches, experiments have been done with light sources that really do emit single photons at wide time intervals, and the same type of interference pattern is observed, showing that a single photon can interfere with *itself*.

The concept of a photon's path is undefined.

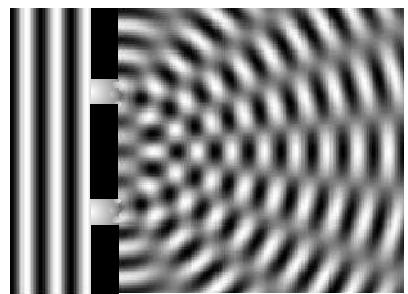
If a single photon can demonstrate double-slit interference, then which slit did it pass through? The unavoidable answer must be that it passes through both! This might not seem so strange if we think of the photon as a wave, but it is highly counterintuitive if we try to visualize it as a particle. The moral is that we should not think in terms of the path of a photon. Like the fully human and fully divine Jesus of Christian theology, a photon is supposed to be 100% wave and 100% particle. If a photon had a well defined path, then it would not demonstrate wave superposition and interference effects, contradicting its wave nature. (In sec. 13.3.4 we will discuss the Heisenberg uncertainty principle, which gives a numerical way of approaching this issue.)

The probability interpretation

The correct interpretation of wave-particle duality is suggested by the random nature of the experiment we've been discussing: even though every photon wave/particle is prepared and released in the same way, the location at which it is eventually detected by the digital camera is different every time. The idea of the probability



j / Bullets pass through a double slit.



k / A water wave passes through a double slit.



l / A single photon can go through both slits.

interpretation of wave-particle duality is that the location of the photon-particle is random, but the probability that it is in a certain location is higher where the photon-wave's amplitude is greater.

More specifically, the probability distribution of the particle must be proportional to the *square* of the wave's amplitude,

$$(\text{probability distribution}) \propto (\text{amplitude})^2.$$

This follows from the correspondence principle and from the fact that a wave's energy density is proportional to the square of its amplitude. If we run the double-slit experiment for a long enough time, the pattern of dots fills in and becomes very smooth as would have been expected in classical physics. To preserve the correspondence between classical and quantum physics, the amount of energy deposited in a given region of the picture over the long run must be proportional to the square of the wave's amplitude. The amount of energy deposited in a certain area depends on the number of photons picked up, which is proportional to the probability of finding any given photon there.

A microwave oven

example 10

► The figure shows two-dimensional (top) and one-dimensional (bottom) representations of the standing wave inside a microwave oven. Gray represents zero field, and white and black signify the strongest fields, with white being a field that is in the opposite direction compared to black. Compare the probabilities of detecting a microwave photon at points A, B, and C.

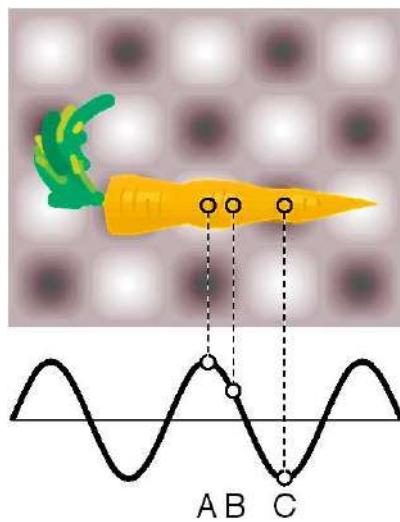
► A and C are both extremes of the wave, so the probabilities of detecting a photon at A and C are equal. It doesn't matter that we have represented C as negative and A as positive, because it is the square of the amplitude that is relevant. The amplitude at B is about 1/2 as much as the others, so the probability of detecting a photon there is about 1/4 as much.

Discussion Questions

A Referring back to the example of the carrot in the microwave oven, show that it would be nonsensical to have probability be proportional to the field itself, rather than the square of the field.

B Einstein did not try to reconcile the wave and particle theories of light, and did not say much about their apparent inconsistency. Einstein basically visualized a beam of light as a stream of bullets coming from a machine gun. In the photoelectric effect, a photon "bullet" would only hit one atom, just as a real bullet would only hit one person. Suppose someone reading his 1905 paper wanted to interpret it by saying that Einstein's so-called particles of light are simply short wave-trains that only occupy a small region of space. Comparing the wavelength of visible light (a few hundred nm) to the size of an atom (on the order of 0.1 nm), explain why this poses a difficulty for reconciling the particle and wave theories.

C Can a white photon exist?



m / Example 10.

D In double-slit diffraction of photons, would you get the same pattern of dots on the digital camera image if you covered one slit? Why should it matter whether you give the photon two choices or only one?

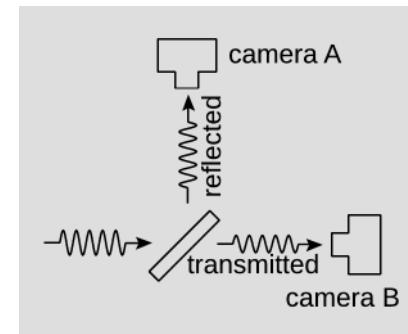
13.2.4 Nonlocality and entanglement

Nonlocality

People sometimes say that quantum mechanics is the set of rules for describing the world of the very small, but this is a false generalization, like saying that terriers are untrainable. How do we define our measure of how small is small? The only distance scales we've discussed have been wavelengths, and there is no upper limit on wavelengths. The wavelength of an FM radio photon is bigger than my terrier, who is very obedient to Newton's laws. The only scale built in to the structure of quantum mechanics is Planck's constant, and Planck's constant has units of joules per hertz, not meters, so it can't be converted into a distance. Quantum mechanics is, as far as we can tell, a valid tool for describing systems at scales from quarks to galaxies.

So quantum behavior can occur at any scale, even large ones. For an example that may be a little disturbing, consider the arrangement shown in figure n. A single photon comes in from the left and encounters a diagonal piece of glass. The glass reflects half the light and transmits half of it. The photon is a wave, and this is expected wave behavior. But the photon is also a particle, and we can't have half a particle. Therefore either camera A will detect a whole photon and B will see none, or it will be the other way around. If we repeat the experiment many times, we might come up with a list of results like this:

A	B
no	yes
yes	no
yes	no
no	yes
no	yes
yes	no
no	yes
yes	no



n / A photon hits a piece of glass that reflects half of the light and transmits the other half.

An instant before the moment of detection, the photon is a wave pattern that just happens to consist of two widely separated pieces, each carrying half the energy. The situation seems perfectly symmetric, but then a moment later we find that B has detected the photon and A hasn't. If B's detection of the photon is random, then how does the information get to A that it had better *not* detect it? This seems as though there is some sort of conspiracy being carried out over arbitrarily large distances and with no time delay. It's as though the two parts of the wave are a pair of criminal suspects who would like to line up their stories but are being kept in separate jail

cells so that they can't communicate. If the part of the wave at B is going to be detected (at full strength, carrying 100% of the energy $E = hf$), how does the part at A get the message that it should fade away like the Cheshire cat? This coordination would have to occur over very large distances — real-world experiments of this type have been done over distances of a thousand kilometers, with the photons traveling either through outer space or through fiber-optic cables. Einstein derisively referred to this apparent coordination as “spooky action at a distance.”

Niels Bohr and two collaborators proposed in 1924 the seemingly reasonable solution that there *can't* be any such coordination. Then the random detection of the photon by camera A and camera B would be independent. Independent probabilities multiply, so there would be a probability of $(1/2)(1/2) = 1/4$ that both cameras would see photons. This would violate conservation of energy, since the original energy $E = hf$ would have been detected twice, and the universe would have gained $1hf$ worth of total energy. But Bohr pointed out that there would also be the same probability that neither camera would detect a photon, in which case the change in the universe's energy would be $-1hf$. On the average, energy would be conserved. According to Bohr's theory, conservation of energy and momentum would not be absolute laws of physics but only rules that would be true on the average.

The experimentalists Geiger and Bothe immediately set out to test this prediction. They performed an experiment analogous to the one in figure n, but with x-rays rather than visible light. Their results, published in 1926, showed that if one detector saw the x-ray photon, the other did not, so that energy was always conserved at the microscopic level, not just on the average. We *never* observe an outcome in which both A and B detect a photon, or one in which neither detects it. That is, the occurrence of event A (camera A sees a photon) and event B (camera B sees one) are both random, but they are not independent.

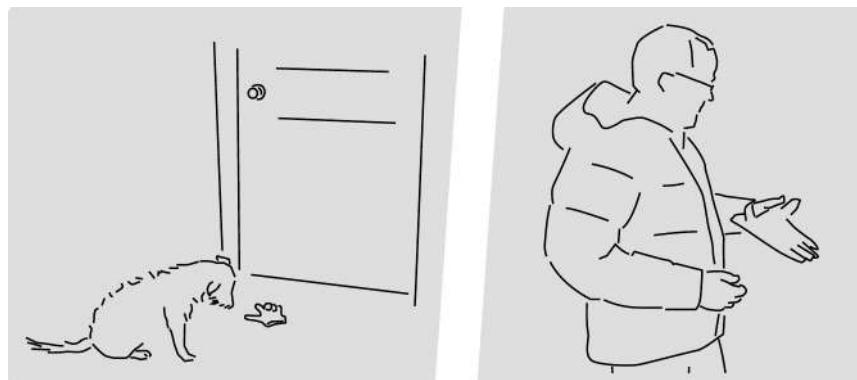
Entanglement

At a 1927 conference in Brussels, Einstein protested that this was a problem, because the two detectors could in principle make their observations simultaneously, and it would then seem that some influence or communication was being transmitted between them faster than the speed of light. “It seems to me,” he complained, “that this difficulty cannot be overcome unless the description of the process in terms of the ... wave is supplemented by some detailed specification of the [trajectory of the particle]. ... If one works only with ... waves, the interpretation ..., I think, contradicts the postulate of relativity.”

The experimental fact ends up being that the spooky action at a distance exists, and it does go faster than light. In 2012, Guerreiro

*et al.*² carried out a very direct and conceptually simple enactment of exactly the experiment in figure n, with electronic timing precise enough to prove that the detection events at A and B were separated from each other by too great a distance to have been linked by any influence traveling at $\leq c$. These findings are summarized by saying that quantum mechanics is *nonlocal*. A single wave-particle can be spread out over an arbitrarily large region of space, but its interactions that transfer energy and momentum are always correlated over these distances in such a way that the conservation laws are maintained.

What Einstein had not originally appreciated was that these correlations do not violate relativity because they do not actually transport any energy, or even any information, between A and B. For example, if Alice is at detector A, and Bob is at B, a million kilometers away, Alice can detect the photon and know immediately that Bob did not detect it. She learns something seemingly instantaneously about Bob — Bob is probably sad and disappointed right now. But because Bob does not have any control over the result, he cannot use this fact to send a message to Alice, so there is no transmission of information. Alice and Bob's states are said to be *entangled*.



o / Entanglement is like finding that you only have your left glove, so that you must have left your right glove at home. There is a gain in information, but no sudden transmission of information from the dog to you.

By analogy, suppose that you head off to work on a winter day in New York. As you step out of the subway station into the cold air, you reach into your pockets for your gloves, but you find that you only have your left glove. Oh, you must have dropped your right glove on the floor while you were petting your adorable terrier on the way out the door. The presence of your left glove tells you that your right glove must be at home. But there has been no spooky action at a distance. You have simply recovered some information about a region of space that lies at some distance from you.

Einstein and Bohr had strong physical intuitions that led them to incorrect predictions about experiments, and these predictions

²arxiv.org/abs/1204.1712. The paper is very readable.

were the fruits of inappropriate mental pictures of what was going on. If we take the principles of quantum mechanics seriously, then the correct picture is the following. Before the photon in figure n hits the glass diagonal, the state of things is the following.

A photon is headed to the right.

Our photon is then partially reflected and partially transmitted. Now we have a superposition of two wave patterns:

c [The photon has been reflected upward.] + c' [The photon has continued to the right.],

where the amplitudes c and c' are equal in absolute value.³

Let's say that the cameras are at equal distances from the glass diagonal, so that their chances to detect the photon occur simultaneously.⁴ After detection, we have this:

c [Camera A detected a photon and B didn't.] + c' [B detected a photon and A didn't.],

Here we have made the nontrivial assumption that material objects like cameras obey the same wave-superposition rules as photons. This turns out to be true. Cameras are made out of things like electrons, and as we'll see in section 13.3, things like electrons are also wave-particles, and they obey all the same wave-particle rules as photons. The states of the two cameras are now entangled.

You can see where this is going. Alice had been standing by camera A, watching anxiously, while Bob, a million kilometers away, was breathlessly observing camera B.

c [Alice saw a photon and Bob didn't. They consider this result to have been random.] + c' [Bob saw a photon and Alice didn't. They consider this result to have been random.],

It doesn't *seem* to Alice and Bob as though their brains are in a superposition of two states. They *feel* as though they have only experienced the one possibility that actually happened, not a mixture of both at the same time. And yet this picture of the physics

³Conservation of energy requires $c^2 = 1/2$ and $c'^2 = 1/2$, even in classical physics. We could have, for example, $c = 1/\sqrt{2}$ and $c' = -1/\sqrt{2}$. Such a possible difference in signs wouldn't concern us in this example. It would only be relevant if there were some later opportunity for the two parts of the wave to recombine and superimpose on one another, producing interference effects.

⁴According to special relativity, this simultaneity holds only in one frame of reference, say the lab frame. But if simultaneity does hold in one frame, then we can also say that in *all* frames, the distance between the two events is "spacelike," i.e., they are too far apart to have been connected by any causal influence propagating at $\leq c$.

explains very nicely how the deterministic laws of physics produce a result that *seems* to them to have been random.

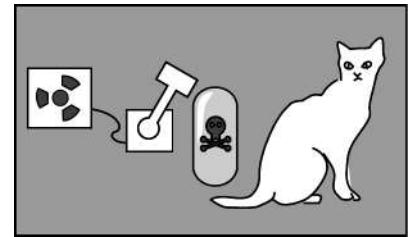
If Alice and Bob have been split into two ghostlike halves of themselves, then conceivably these half-selves could undergo interference, as in the double-slit experiment. But there are practical reasons why we cannot actually detect such interference effects. For one thing, Alice and Bob are macroscopic objects, with energies E on the order of many joules. Because Planck's constant is small, their wave frequencies $f = E/h$ are extremely high, and their wavelengths incredibly short (on the order of 10^{-34} m!). We have seen that diffraction becomes undetectable when wavelengths are too short. Furthermore, there is a phenomenon called decoherence, discussed further in sec. 14.9.2, p. 1002, in which interactions with the environment tend to rapidly randomize the wave-phases of large objects. When phases are randomized, interference and diffraction effects become undetectable.

Historically, it seemed absurd to the originators of quantum mechanics to imagine a macroscopic object in a superposition of states. The most celebrated example is called the Schrödinger's cat experiment. Luckily for the cat, there probably was no actual experiment — it was simply a “thought experiment” that the German theorist Schrödinger discussed with his colleagues. Schrödinger wrote:

One can even construct quite burlesque cases. A cat is shut up in a steel container, together with the following diabolical apparatus (which one must keep out of the direct clutches of the cat): In a Geiger tube [radiation detector] there is a tiny mass of radioactive substance, so little that in the course of an hour perhaps one atom of it disintegrates, but also with equal probability not even one; if it does happen, the counter [detector] responds and ... activates a hammer that shatters a little flask of prussic acid [filling the chamber with poison gas]. If one has left this entire system to itself for an hour, then one will say to himself that the cat is still living, if in that time no atom has disintegrated. The first atomic disintegration would have poisoned it.

It seemed ridiculous to Schrödinger that at the end of the hour, “The uncertainty originally restricted to the atomic domain has been transformed into a macroscopic uncertainty...,” and the cat would be in a superposed state.

In modern language, people like Einstein and Schrödinger didn't feel comfortable with nonlocality, or with entanglement of subatomic particles, and they felt even less comfortable with applying these concepts to macroscopic objects. Today, entanglement has been



p / Schrödinger's cat.

demonstrated using objects that clearly deserve to be called macroscopic. For example, in 2012, K.C. Lee *et al.* created a version of the experiment in figure n in which the cameras were replaced by small diamonds, about 1 mm in size. They were separated by 15 cm, which is a macroscopic distance. When a photon hit one of the diamonds, it produced a vibration in the crystal lattice. This vibration was localized to a relatively small region within the diamond, but this region was still large enough that one has to admit that it qualifies as macroscopic. Its atoms had a total weight of about 0.1 nanograms, which is a quantity big enough to weigh on a state-of-the-art balance, and the region was about 0.01 mm in size, which would make it visible with a magnifying glass.

The quantum states of the two diamonds became entangled: if one had detected the photon, the other hadn't. This entangled state was maintained for only about 7 picoseconds before decoherence destroyed the phase relationship between one diamond and the other. But Lee was able to use additional photons to “read out” the quantum states in only 0.5 ps, before decoherence occurred, and verify that there were wave interference effects in which one diamond's quantum-mechanical wave had a definite phase relationship with the other's. Although these experiments are difficult, they suggest that there is no obstruction in principle to observing quantum-mechanical effects such as superposition in arbitrarily large objects.

Entanglement is discussed in more mathematical detail in sec. 14.11, p. 1007.

The Copenhagen and many-worlds approximations

When we last saw Alice and Bob, they were in this superposition of states,

$$c \begin{cases} \text{Alice saw a photon} \\ \text{and Bob didn't. They} \\ \text{consider this result to} \\ \text{have been random.} \end{cases} + c' \begin{cases} \text{Bob saw a photon and} \\ \text{Alice didn't. They} \\ \text{consider this result to} \\ \text{have been random.} \end{cases}$$

with

$$|c| = |c'|.$$

Let's focus on Bob number one — the sad Bob — who didn't see a photon. This is just one of the disappointments that Bob has experienced in his life, which include breaking up with his college crush and failing to summit Kilimanjaro due to altitude sickness. But Bob is a sane, normal person, and he's not going to spend the rest of his life obsessing over how things might have been, in another world. Like a banker writing off a bad debt, Bob decides to stop maintaining all the bookkeeping that is, to him, irrelevant going forward. He now rewrites history and says that

$$|c| = 1 \quad \text{and} \quad |c'| = 0.$$

Technically speaking, this is wrong, because it is in principle still possible to have wave interference effects between sad Bob and happy Bob. But such effects are impractical to observe, due to effects like short wavelengths and decoherence, so what Bob is doing by “clearing the books” is an extremely good approximation. We will refer to this approximation by two different names, the *Copenhagen approximation* and the *many-worlds approximation*, for the following historical and psychological reasons.

In the early years of quantum mechanics, the school of physicists centering on Niels Bohr in Copenhagen were horribly confused about how to interpret quantum mechanics. They had all kinds of wrong ideas, such as the idea that quantum mechanics applied to individual atoms but not to light or to macroscopic objects. They didn’t know about decoherence. They thought there was a clear dividing line between microscopic things and macroscopic things (there isn’t), and they hypothesized that quantum mechanics only applied to microscopic ones (it applies to both). They claimed that clearing the books was an actual physical process, which they described as the “collapse” of the wave. This has traditionally been referred to as the Copenhagen “interpretation,” but we can now see that it is an approximation. There are cases where it is a bad approximation, e.g., at $t = 3$ ps during the experiment by Lee *et al.* (p. 888), when decoherence had started to happen but was only about half-way complete.

The many-worlds approximation came along a little later.⁵ It consists of making the same “clearing the books” approximation, but recognizing that there is no physical process of collapse.

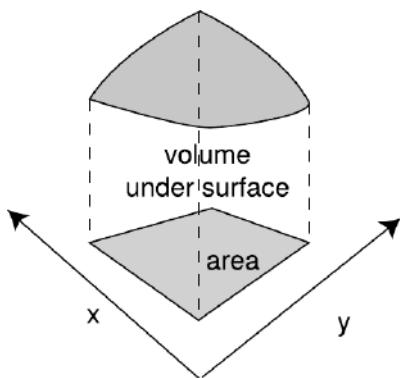
Many physicists are philosophically attached to one or the other of these approximations, and would object to my description of them as approximations. My main purpose in writing this explanation is to immunize you against the impression, which can be mistakenly picked up from many descriptions of quantum mechanics, that a particular point of view on these topics (often the Copenhagen approximation) is somehow “standard.”

13.2.5 Photons in three dimensions

Up until now I’ve been sneaky and avoided a full discussion of the three-dimensional aspects of the probability interpretation. The example of the carrot in the microwave oven, for example, reduced to a one-dimensional situation because we were considering three points along the same line and because we were only comparing ratios of probabilities.

A typical example of a probability distribution in section 13.1

⁵It was originally proposed in a 1957 PhD thesis by Hugh Everett, who called it the relative state interpretation of quantum mechanics. Later it began to be referred to as the many-worlds interpretation.



q / Probability is the volume under a surface defined by $D(x, y)$.

was the distribution of heights of human beings. The thing that varied randomly, height, h , had units of meters, and the probability distribution was a graph of a function $D(h)$. The units of the probability distribution had to be m^{-1} (inverse meters) so that areas under the curve, interpreted as probabilities, would be unitless: $(\text{area}) = (\text{height})(\text{width}) = \text{m}^{-1} \cdot \text{m}$.

Now suppose we have a two-dimensional problem, e.g., the probability distribution for the place on the surface of a digital camera chip where a photon will be detected. The point where it is detected would be described with two variables, x and y , each having units of meters. The probability distribution will be a function of both variables, $D(x, y)$. A probability is now visualized as the volume under the surface described by the function $D(x, y)$, as shown in figure q. The units of D must be m^{-2} so that probabilities will be unitless: $(\text{probability}) = (\text{depth})(\text{length})(\text{width}) = \text{m}^{-2} \cdot \text{m} \cdot \text{m}$. In terms of calculus, we have $P = \int D \, dx \, dy$.

Generalizing finally to three dimensions, we find by analogy that the probability distribution will be a function of all three coordinates, $D(x, y, z)$, and will have units of m^{-3} . It is unfortunately impossible to visualize the graph unless you are a mutant with a natural feel for life in four dimensions. If the probability distribution is nearly constant within a certain volume of space v , the probability that the photon is in that volume is simply vD . If not, then we can use an integral, $P = \int D \, dx \, dy \, dz$.



13.3 Matter as a wave

[In] a few minutes I shall be all melted... I have been wicked in my day, but I never thought a little girl like you would ever be able to melt me and end my wicked deeds. Look out — here I go!

The Wicked Witch of the West

As the Wicked Witch learned the hard way, losing molecular cohesion can be unpleasant. That's why we should be very grateful that the concepts of quantum physics apply to matter as well as light. If matter obeyed the laws of classical physics, molecules wouldn't exist.

Consider, for example, the simplest atom, hydrogen. Why does one hydrogen atom form a chemical bond with another hydrogen atom? Roughly speaking, we'd expect a neighboring pair of hydrogen atoms, A and B, to exert no force on each other at all, attractive or repulsive: there are two repulsive interactions (proton A with proton B and electron A with electron B) and two attractive interactions (proton A with electron B and electron A with proton B). Thinking a little more precisely, we should even expect that once the two atoms got close enough, the interaction would be repulsive. For instance, if you squeezed them so close together that the two protons were almost on top of each other, there would be a tremendously strong repulsion between them due to the $1/r^2$ nature of the electrical force. A more detailed calculation using classical physics gives an extremely weak binding, about $1/17$ the strength of what we actually observe, which is far too weak to make the bond hold together.

Quantum physics to the rescue! As we'll see shortly, the whole problem is solved by applying the same quantum concepts to elec-

trons that we have already used for photons.

13.3.1 Electrons as waves

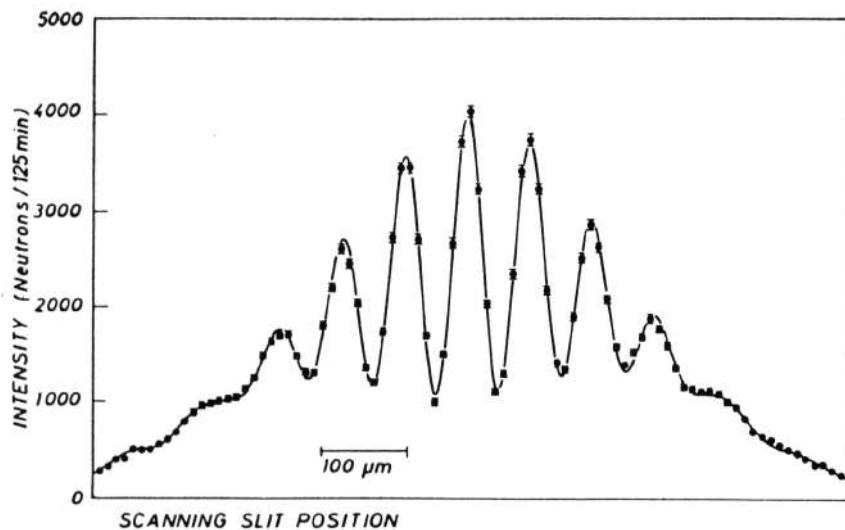
We started our journey into quantum physics by studying the random behavior of *matter* in radioactive decay, and then asked how randomness could be linked to the basic laws of nature governing *light*. The probability interpretation of wave-particle duality was strange and hard to accept, but it provided such a link. It is now natural to ask whether the same explanation could be applied to matter. If the fundamental building block of light, the photon, is a particle as well as a wave, is it possible that the basic units of matter, such as electrons, are waves as well as particles?

A young French aristocrat studying physics, Louis de Broglie (pronounced “broylee”), made exactly this suggestion in his 1923 Ph.D. thesis. His idea had seemed so farfetched that there was serious doubt about whether to grant him the degree. Einstein was asked for his opinion, and with his strong support, de Broglie got his degree.

Only two years later, American physicists C.J. Davisson and L. Germer confirmed de Broglie’s idea by accident. They had been studying the scattering of electrons from the surface of a sample of nickel, made of many small crystals. (One can often see such a crystalline pattern on a brass doorknob that has been polished by repeated handling.) An accidental explosion occurred, and when they put their apparatus back together they observed something entirely different: the scattered electrons were now creating an interference pattern! This dramatic proof of the wave nature of matter came about because the nickel sample had been melted by the explosion and then resolidified as a single crystal. The nickel atoms, now nicely arranged in the regular rows and columns of a crystalline lattice, were acting as the lines of a diffraction grating. The new crystal was analogous to the type of ordinary diffraction grating in which the lines are etched on the surface of a mirror (a reflection grating) rather than the kind in which the light passes through the transparent gaps between the lines (a transmission grating).

Although we will concentrate on the wave-particle duality of electrons because it is important in chemistry and the physics of atoms, all the other “particles” of matter you’ve learned about show wave properties as well. Figure a, for instance, shows a wave interference pattern of neutrons.

It might seem as though all our work was already done for us, and there would be nothing new to understand about electrons: they have the same kind of funny wave-particle duality as photons. That’s almost true, but not quite. There are some important ways in which electrons differ significantly from photons:



a / A double-slit interference pattern made with neutrons. (A. Zeilinger, R. Gähler, C.G. Shull, W. Treimer, and W. Mampe, *Reviews of Modern Physics*, Vol. 60, 1988.)

1. Electrons have mass, and photons don't.
2. Photons always move at the speed of light, but electrons can move at any speed less than c .
3. Photons don't have electric charge, but electrons do, so electric forces can act on them. The most important example is the atom, in which the electrons are held by the electric force of the nucleus.
4. Electrons cannot be absorbed or emitted as photons are. Destroying an electron or creating one out of nothing would violate conservation of charge.

(In section 13.4 we will learn of one more fundamental way in which electrons differ from photons, for a total of five.)

Because electrons are different from photons, it is not immediately obvious which of the photon equations from chapter 11 can be applied to electrons as well. A particle property, the energy of one photon, is related to its wave properties via $E = hf$ or, equivalently, $E = hc/\lambda$. The momentum of a photon was given by $p = hf/c$ or $p = h/\lambda$. Ultimately it was a matter of experiment to determine which of these equations, if any, would work for electrons, but we can make a quick and dirty guess simply by noting that some of the equations involve c , the speed of light, and some do not. Since c is irrelevant in the case of an electron, we might guess that the

equations of general validity are those that do not have c in them:

$$E = hf$$
$$p = h/\lambda$$

This is essentially the reasoning that de Broglie went through, and experiments have confirmed these two equations for all the fundamental building blocks of light and matter, not just for photons and electrons.

The second equation, which I soft-pedaled in the previous chapter, takes on a greater importance for electrons. This is first of all because the momentum of matter is more likely to be significant than the momentum of light under ordinary conditions, and also because force is the transfer of momentum, and electrons are affected by electrical forces.

The wavelength of an elephant

example 11

▷ What is the wavelength of a trotting elephant?

▷ One may doubt whether the equation should be applied to an elephant, which is not just a single particle but a rather large collection of them. Throwing caution to the wind, however, we estimate the elephant's mass at 10^3 kg and its trotting speed at 10 m/s. Its wavelength is therefore roughly

$$\begin{aligned}\lambda &= \frac{h}{p} \\ &= \frac{h}{mv} \\ &= \frac{6.63 \times 10^{-34} \text{ J}\cdot\text{s}}{(10^3 \text{ kg})(10 \text{ m/s})} \\ &\sim 10^{-37} \frac{(\text{kg}\cdot\text{m}^2/\text{s}^2)\cdot\text{s}}{\text{kg}\cdot\text{m/s}} \\ &= 10^{-37} \text{ m}\end{aligned}$$

The wavelength found in this example is so fantastically small that we can be sure we will never observe any measurable wave phenomena with elephants or any other human-scale objects. The result is numerically small because Planck's constant is so small, and as in some examples encountered previously, this smallness is in accord with the correspondence principle.

Although a smaller mass in the equation $\lambda = h/mv$ does result in a longer wavelength, the wavelength is still quite short even for individual electrons under typical conditions, as shown in the following example.

The typical wavelength of an electron

example 12

▷ Electrons in circuits and in atoms are typically moving through

voltage differences on the order of 1 V, so that a typical energy is $(e)(1 \text{ V})$, which is on the order of 10^{-19} J . What is the wavelength of an electron with this amount of kinetic energy?

▷ This energy is nonrelativistic, since it is much less than mc^2 . Momentum and energy are therefore related by the nonrelativistic equation $K = p^2/2m$. Solving for p and substituting in to the equation for the wavelength, we find

$$\begin{aligned}\lambda &= \frac{h}{\sqrt{2mK}} \\ &= 1.6 \times 10^{-9} \text{ m.}\end{aligned}$$

This is on the same order of magnitude as the size of an atom, which is no accident: as we will discuss in the next chapter in more detail, an electron in an atom can be interpreted as a standing wave. The smallness of the wavelength of a typical electron also helps to explain why the wave nature of electrons wasn't discovered until a hundred years after the wave nature of light. To scale the usual wave-optics devices such as diffraction gratings down to the size needed to work with electrons at ordinary energies, we need to make them so small that their parts are comparable in size to individual atoms. This is essentially what Davisson and Germer did with their nickel crystal.

self-check E

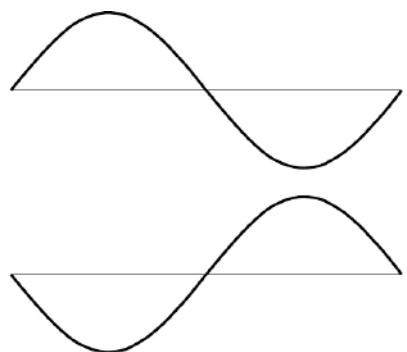
These remarks about the inconvenient smallness of electron wavelengths apply only under the assumption that the electrons have typical energies. What kind of energy would an electron have to have in order to have a longer wavelength that might be more convenient to work with?

▷ Answer, p. 1067

What kind of wave is it?

If a sound wave is a vibration of matter, and a photon is a vibration of electric and magnetic fields, what kind of a wave is an electron made of? The disconcerting answer is that there is no experimental “observable,” i.e., directly measurable quantity, to correspond to the electron wave itself. In other words, there are devices like microphones that detect the oscillations of air pressure in a sound wave, and devices such as radio receivers that measure the oscillation of the electric and magnetic fields in a light wave, but nobody has ever found any way to measure the electron wave directly.

We can of course detect the energy (or momentum) possessed by an electron just as we could detect the energy of a photon using a digital camera. (In fact I'd imagine that an unmodified digital camera chip placed in a vacuum chamber would detect electrons just as handily as photons.) But this only allows us to determine where the



b / These two electron waves are not distinguishable by any measuring device.

wave carries high probability and where it carries low probability. Probability is proportional to the square of the wave's amplitude, but measuring its square is not the same as measuring the wave itself. In particular, we get the same result by squaring either a positive number or its negative, so there is no way to determine the positive or negative sign of an electron wave. This unobservability of the phase of the wavefunction is discussed in more detail on p. 917.

Most physicists tend toward the school of philosophy known as operationalism, which says that a concept is only meaningful if we can define some set of operations for observing, measuring, or testing it. According to a strict operationalist, then, the electron wave itself is a meaningless concept. Nevertheless, it turns out to be one of those concepts like love or humor that is impossible to measure and yet very useful to have around. We therefore give it a symbol, Ψ (the capital Greek letter psi), and a special name, the electron *wavefunction* (because it is a function of the coordinates x , y , and z that specify where you are in space). It would be impossible, for example, to calculate the shape of the electron wave in a hydrogen atom without having some symbol for the wave. But when the calculation produces a result that can be compared directly to experiment, the final algebraic result will turn out to involve only Ψ^2 , which is what is observable, not Ψ itself.

Since Ψ , unlike E and B , is not directly measurable, we are free to make the probability equations have a simple form: instead of having the probability density equal to some funny constant multiplied by Ψ^2 , we simply define Ψ so that the constant of proportionality is one:

$$(\text{probability distribution}) = |\Psi|^2.$$

Since the probability distribution has units of m^{-3} , the units of Ψ must be $\text{m}^{-3/2}$. The square of a negative number is still positive, so the absolute value signs may seem unnecessary, but as we'll see on p. 913 in sec. 13.3.6, the wavefunction may in general be a complex number. In fact, only standing waves, not traveling waves, can really be represented by real numbers, although we will often cheat and draw pictures of traveling waves as if they were real-valued functions.

Discussion Question

A Frequency is oscillations per second, whereas wavelength is meters per oscillation. How could the equations $E = hf$ and $p = h/\lambda$ be made to look more alike by using quantities that were more closely analogous? (This more symmetric treatment makes it easier to incorporate relativity into quantum mechanics, since relativity says that space and time are not entirely separate.)

13.3.2 Dispersive waves

A colleague of mine who teaches chemistry loves to tell the story about an exceptionally bright student who, when told of the equation $p = h/\lambda$, protested, "But when I derived it, it had a factor of

2!" The issue that's involved is a real one, albeit one that could be glossed over (and is, in most textbooks) without raising any alarms in the mind of the average student. The present optional section addresses this point; it is intended for the student who wishes to delve a little deeper.

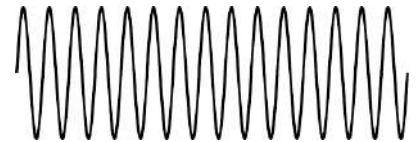
Here's how the now-legendary student was presumably reasoning. We start with the equation $v = f\lambda$, which is valid for any sine wave, whether it's quantum or classical. Let's assume we already know $E = hf$, and are trying to derive the relationship between wavelength and momentum:

$$\begin{aligned}\lambda &= \frac{v}{f} \\ &= \frac{vh}{E} \\ &= \frac{vh}{\frac{1}{2}mv^2} \\ &= \frac{2h}{mv} \\ &= \frac{2h}{p}.\end{aligned}$$

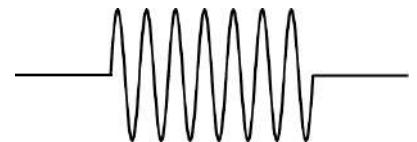
The reasoning seems valid, but the result does contradict the accepted one, which is after all solidly based on experiment.

The mistaken assumption is that we can figure everything out in terms of pure sine waves. Mathematically, the only wave that has a perfectly well defined wavelength and frequency is a sine wave, and not just any sine wave but an infinitely long sine wave, c. The unphysical thing about such a wave is that it has no leading or trailing edge, so it can never be said to enter or leave any particular region of space. Our derivation made use of the velocity, v , and if velocity is to be a meaningful concept, it must tell us how quickly stuff (mass, energy, momentum, ...) is transported from one region of space to another. Since an infinitely long sine wave doesn't remove any stuff from one region and take it to another, the "velocity of its stuff" is not a well defined concept.

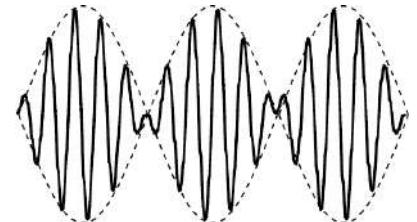
Of course the individual wave peaks do travel through space, and one might think that it would make sense to associate their speed with the "speed of stuff," but as we will see, the two velocities are in general unequal when a wave's velocity depends on wavelength. Such a wave is called a *dispersive* wave, because a wave pulse consisting of a superposition of waves of different wavelengths will separate (disperse) into its separate wavelengths as the waves move through space at different speeds. Nearly all the waves we have encountered have been nondispersive. For instance, sound waves and light waves (in a vacuum) have speeds independent of wavelength. A water wave is one good example of a dispersive wave. Long-wavelength water



c / Part of an infinite sine wave.



d / A finite-length sine wave.



e / A beat pattern created by superimposing two sine waves with slightly different wavelengths.

waves travel faster, so a ship at sea that encounters a storm typically sees the long-wavelength parts of the wave first. When dealing with dispersive waves, we need symbols and words to distinguish the two speeds. The speed at which wave peaks move is called the phase velocity, v_p , and the speed at which “stuff” moves is called the group velocity, v_g .

An infinite sine wave can only tell us about the phase velocity, not the group velocity, which is really what we would be talking about when we refer to the speed of an electron. If an infinite sine wave is the simplest possible wave, what’s the next best thing? We might think the runner up in simplicity would be a wave train consisting of a chopped-off segment of a sine wave, d. However, this kind of wave has kinks in it at the end. A simple wave should be one that we can build by superposing a small number of infinite sine waves, but a kink can never be produced by superposing any number of infinitely long sine waves.

Actually the simplest wave that transports stuff from place to place is the pattern shown in figure e. Called a beat pattern, it is formed by superposing two sine waves whose wavelengths are similar but not quite the same. If you have ever heard the pulsating howling sound of musicians in the process of tuning their instruments to each other, you have heard a beat pattern. The beat pattern gets stronger and weaker as the two sine waves go in and out of phase with each other. The beat pattern has more “stuff” (energy, for example) in the areas where constructive interference occurs, and less in the regions of cancellation. As the whole pattern moves through space, stuff is transported from some regions and into other ones.

If the frequency of the two sine waves differs by 10%, for instance, then ten periods will occur between times when they are in phase. Another way of saying it is that the sinusoidal “envelope” (the dashed lines in figure e) has a frequency equal to the difference in frequency between the two waves. For instance, if the waves had frequencies of 100 Hz and 110 Hz, the frequency of the envelope would be 10 Hz.

To apply similar reasoning to the wavelength, we must define a quantity $z = 1/\lambda$ that relates to wavelength in the same way that frequency relates to period. In terms of this new variable, the z of the envelope equals the difference between the z' s of the two sine waves.

The group velocity is the speed at which the envelope moves through space. Let Δf and Δz be the differences between the frequencies and z' s of the two sine waves, which means that they equal the frequency and z of the envelope. The group velocity is $v_g = f_{\text{envelope}}\lambda_{\text{envelope}} = \Delta f/\Delta z$. If Δf and Δz are sufficiently

small, we can approximate this expression as a derivative,

$$v_g = \frac{df}{dz}.$$

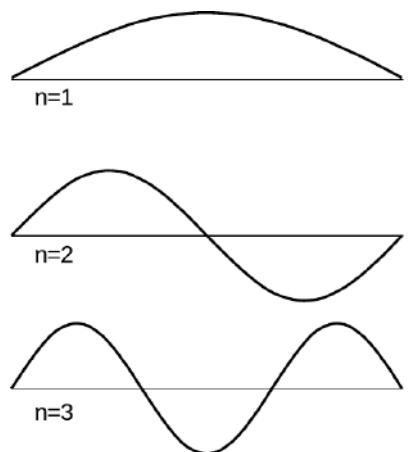
This expression is usually taken as the definition of the group velocity for wave patterns that consist of a superposition of sine waves having a narrow range of frequencies and wavelengths. In quantum mechanics, with $f = E/h$ and $z = p/h$, we have $v_g = dE/dp$. In the case of a nonrelativistic electron the relationship between energy and momentum is $E = p^2/2m$, so the group velocity is $dE/dp = p/m = v$, exactly what it should be. It is only the phase velocity that differs by a factor of two from what we would have expected, but the phase velocity is not the physically important thing.

13.3.3 Bound states

Electrons are at their most interesting when they're in atoms, that is, when they are bound within a small region of space. We can understand a great deal about atoms and molecules based on simple arguments about such bound states, without going into any of the realistic details of atom. The simplest model of a bound state is known as the particle in a box: like a ball on a pool table, the electron feels zero force while in the interior, but when it reaches an edge it encounters a wall that pushes back inward on it with a large force. In particle language, we would describe the electron as bouncing off of the wall, but this incorrectly assumes that the electron has a certain path through space. It is more correct to describe the electron as a wave that undergoes 100% reflection at the boundaries of the box.

Like generations of physics students before me, I rolled my eyes when initially introduced to the unrealistic idea of putting a particle in a box. It seemed completely impractical, an artificial textbook invention. Today, however, it has become routine to study electrons in rectangular boxes in actual laboratory experiments. The “box” is actually just an empty cavity within a solid piece of silicon, amounting in volume to a few hundred atoms. The methods for creating these electron-in-a-box setups (known as “quantum dots”) were a by-product of the development of technologies for fabricating computer chips.

For simplicity let's imagine a one-dimensional electron in a box, i.e., we assume that the electron is only free to move along a line. The resulting standing wave patterns, of which the first three are shown in the figure, are just like some of the patterns we encountered with sound waves in musical instruments. The wave patterns must be zero at the ends of the box, because we are assuming the walls are impenetrable, and there should therefore be zero probability of finding the electron outside the box. Each wave pattern is labeled according to n , the number of peaks and valleys it has. In



f / Three possible standing-wave patterns for a particle in a box.

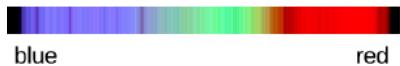
quantum physics, these wave patterns are referred to as “states” of the particle-in-the-box system.

The following seemingly innocuous observations about the particle in the box lead us directly to the solutions to some of the most vexing failures of classical physics:

The particle’s energy is quantized (can only have certain values). Each wavelength corresponds to a certain momentum, and a given momentum implies a definite kinetic energy, $E = p^2/2m$. (This is the second type of energy quantization we have encountered. The type we studied previously had to do with restricting the number of particles to a whole number, while assuming some specific wavelength and energy for each particle. This type of quantization refers to the energies that a single particle can have. Both photons and matter particles demonstrate both types of quantization under the appropriate circumstances.)

The particle has a minimum kinetic energy. Long wavelengths correspond to low momenta and low energies. There can be no state with an energy lower than that of the $n = 1$ state, called the ground state.

The smaller the space in which the particle is confined, the higher its kinetic energy must be. Again, this is because long wavelengths give lower energies.



g / The spectrum of the light from the star Sirius.

Spectra of thin gases

example 13

A fact that was inexplicable by classical physics was that thin gases absorb and emit light only at certain wavelengths. This was observed both in earthbound laboratories and in the spectra of stars. The figure on the left shows the example of the spectrum of the star Sirius, in which there are “gap teeth” at certain wavelengths. Taking this spectrum as an example, we can give a straightforward explanation using quantum physics.

Energy is released in the dense interior of the star, but the outer layers of the star are thin, so the atoms are far apart and electrons are confined within individual atoms. Although their standing-wave patterns are not as simple as those of the particle in the box, their energies are quantized.

When a photon is on its way out through the outer layers, it can be absorbed by an electron in an atom, but only if the amount of energy it carries happens to be the right amount to kick the electron from one of the allowed energy levels to one of the higher levels. The photon energies that are missing from the spectrum are the ones that equal the difference in energy between two electron energy levels. (The most prominent of the absorption lines in Sirius’s spectrum are absorption lines of the hydrogen atom.)

The stability of atoms

example 14

In many Star Trek episodes the Enterprise, in orbit around a planet, suddenly lost engine power and began spiraling down toward the planet's surface. This was utter nonsense, of course, due to conservation of energy: the ship had no way of getting rid of energy, so it did not need the engines to replenish it.

Consider, however, the electron in an atom as it orbits the nucleus. The electron *does* have a way to release energy: it has an acceleration due to its continuously changing direction of motion, and according to classical physics, any accelerating charged particle emits electromagnetic waves. According to classical physics, atoms should collapse!

The solution lies in the observation that a bound state has a minimum energy. An electron in one of the higher-energy atomic states can and does emit photons and hop down step by step in energy. But once it is in the ground state, it cannot emit a photon because there is no lower-energy state for it to go to.

Chemical bonds

example 15

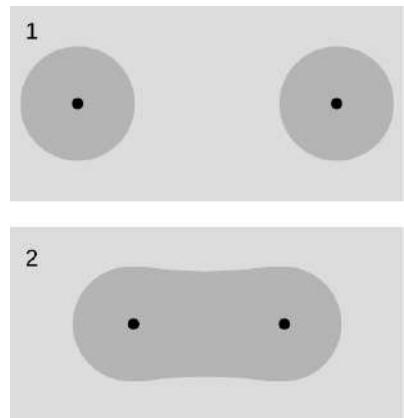
I began this section with a classical argument that chemical bonds, as in an H_2 molecule, should not exist. Quantum physics explains why this type of bonding does in fact occur. There are actually two effects going on, one due to kinetic energy and one due to electrical energy. We'll concentrate on the kinetic energy effect in this example. Example 24 on page 934 revisits the H_2 bond in more detail. (A qualitatively different type of bonding is discussed on page 941.)

The kinetic energy effect is pretty simple. When the atoms are next to each other, the electrons are shared between them. The "box" is about twice as wide, and a larger box allows a smaller kinetic energy. Energy is required in order to separate the atoms.

Discussion Questions

A Neutrons attract each other via the strong nuclear force, so according to classical physics it should be possible to form nuclei out of clusters of two or more neutrons, with no protons at all. Experimental searches, however, have failed to turn up evidence of a stable two-neutron system (dineutron) or larger stable clusters. These systems are apparently not just unstable in the sense of being able to beta decay but unstable in the sense that they don't hold together at all. Explain based on quantum physics why a dineutron might spontaneously fly apart.

B The following table shows the energy gap between the ground state and the first excited state for four nuclei, in units of picojoules. (The nuclei were chosen to be ones that have similar structures, e.g., they are all spherical in shape.)



h / Two hydrogen atoms bond to form an H_2 molecule. In the molecule, the two electrons' wave patterns overlap , and are about twice as wide.

nucleus	energy gap (picojoules)
^4He	3.234
^{16}O	0.968
^{40}Ca	0.536
^{208}Pb	0.418

Explain the trend in the data.

13.3.4 The uncertainty principle

Eliminating randomness through measurement?

A common reaction to quantum physics, among both early-twentieth-century physicists and modern students, is that we should be able to get rid of randomness through accurate measurement. If I say, for example, that it is meaningless to discuss the path of a photon or an electron, one might suggest that we simply measure the particle's position and velocity many times in a row. This series of snapshots would amount to a description of its path.

A practical objection to this plan is that the process of measurement will have an effect on the thing we are trying to measure. This may not be of much concern, for example, when a traffic cop measures your car's motion with a radar gun, because the energy and momentum of the radar pulses are insufficient to change the car's motion significantly. But on the subatomic scale it is a very real problem. Making a videotape through a microscope of an electron orbiting a nucleus is not just difficult, it is theoretically impossible. The video camera makes pictures of things using light that has bounced off them and come into the camera. If even a single photon of visible light was to bounce off of the electron we were trying to study, the electron's recoil would be enough to change its behavior significantly.

The Heisenberg uncertainty principle

This insight, that measurement changes the thing being measured, is the kind of idea that clove-cigarette-smoking intellectuals outside of the physical sciences like to claim they knew all along. If only, they say, the physicists had made more of a habit of reading literary journals, they could have saved a lot of work. The anthropologist Margaret Mead has recently been accused of inadvertently encouraging her teenaged Samoan informants to exaggerate the freedom of youthful sexual experimentation in their society. If this is considered a damning critique of her work, it is because she could have done better: other anthropologists claim to have been able to eliminate the observer-as-participant problem and collect untainted data.

The German physicist Werner Heisenberg, however, showed that in quantum physics, *any* measuring technique runs into a brick wall when we try to improve its accuracy beyond a certain point. Heisenberg showed that the limitation is a question of *what there is to be*



i / Werner Heisenberg (1901–1976). Heisenberg helped to develop the foundations of quantum mechanics, including the Heisenberg uncertainty principle. He was the scientific leader of the Nazi atomic-bomb program up until its cancellation in 1942, when the military decided that it was too ambitious a project to undertake in wartime, and too unlikely to produce results.

known, even in principle, about the system itself, not of the ability or inability of a specific measuring device to ferret out information that is knowable but not previously hidden.

Suppose, for example, that we have constructed an electron in a box (quantum dot) setup in our laboratory, and we are able to adjust the length L of the box as desired. All the standing wave patterns pretty much fill the box, so our knowledge of the electron's position is of limited accuracy. If we write Δx for the range of uncertainty in our knowledge of its position, then Δx is roughly the same as the length of the box:

$$\Delta x \approx L$$

If we wish to know its position more accurately, we can certainly squeeze it into a smaller space by reducing L , but this has an unintended side-effect. A standing wave is really a superposition of two traveling waves going in opposite directions. The equation $p = h/\lambda$ really only gives the magnitude of the momentum vector, not its direction, so we should really interpret the wave as a 50/50 mixture of a right-going wave with momentum $p = h/\lambda$ and a left-going one with momentum $p = -h/\lambda$. The uncertainty in our knowledge of the electron's momentum is $\Delta p = 2h/\lambda$, covering the range between these two values. Even if we make sure the electron is in the ground state, whose wavelength $\lambda = 2L$ is the longest possible, we have an uncertainty in momentum of $\Delta p = h/L$. In general, we find

$$\Delta p \gtrsim h/L,$$

with equality for the ground state and inequality for the higher-energy states. Thus if we reduce L to improve our knowledge of the electron's position, we do so at the cost of knowing less about its momentum. This trade-off is neatly summarized by multiplying the two equations to give

$$\Delta p \Delta x \gtrsim h.$$

Although we have derived this in the special case of a particle in a box, it is an example of a principle of more general validity:

The Heisenberg uncertainty principle

It is not possible, even in principle, to know the momentum and the position of a particle simultaneously and with perfect accuracy. The uncertainties in these two quantities are always such that $\Delta p \Delta x \gtrsim h$.

(This approximation can be made into a strict inequality, $\Delta p \Delta x > h/4\pi$, but only with more careful definitions, which we will not bother with.)

Note that although I encouraged you to think of this derivation in terms of a specific real-world system, the quantum dot, no

reference was ever made to any specific laboratory equipment or procedures. The argument is simply that we cannot *know* the particle's position very accurately unless it *has* a very well defined position, it cannot have a very well defined position unless its wave-pattern covers only a very small amount of space, and its wave-pattern cannot be thus compressed without giving it a short wavelength and a correspondingly uncertain momentum. The uncertainty principle is therefore a restriction on how much there is to know about a particle, not just on what we can know about it with a certain technique.

An estimate for electrons in atoms

example 16

- ▷ A typical energy for an electron in an atom is on the order of $(1 \text{ volt}) \cdot e$, which corresponds to a speed of about 1% of the speed of light. If a typical atom has a size on the order of 0.1 nm, how close are the electrons to the limit imposed by the uncertainty principle?
- ▷ If we assume the electron moves in all directions with equal probability, the uncertainty in its momentum is roughly twice its typical momentum. This is only an order-of-magnitude estimate, so we take Δp to be the same as a typical momentum:

$$\begin{aligned}\Delta p \Delta x &= p_{\text{typical}} \Delta x \\ &= (m_{\text{electron}})(0.01c)(0.1 \times 10^{-9} \text{ m}) \\ &= 3 \times 10^{-34} \text{ J}\cdot\text{s}\end{aligned}$$

This is on the same order of magnitude as Planck's constant, so evidently the electron is "right up against the wall." (The fact that it is somewhat less than h is of no concern since this was only an estimate, and we have not stated the uncertainty principle in its most exact form.)

self-check F

If we were to apply the uncertainty principle to human-scale objects, what would be the significance of the small numerical value of Planck's constant?

▷ Answer, p. 1067

Discussion Questions

- A** Compare Δp and Δx for the two lowest energy levels of the one-dimensional particle in a box, and discuss how this relates to the uncertainty principle.
- B** On a graph of Δp versus Δx , sketch the regions that are allowed and forbidden by the Heisenberg uncertainty principle. Interpret the graph: Where does an atom lie on it? An elephant? Can either p or x be measured with perfect accuracy if we don't care about the other?

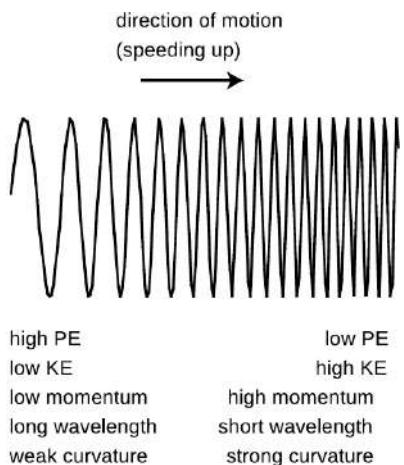
13.3.5 Electrons in electric fields

So far the only electron wave patterns we've considered have been simple sine waves, but whenever an electron finds itself in an

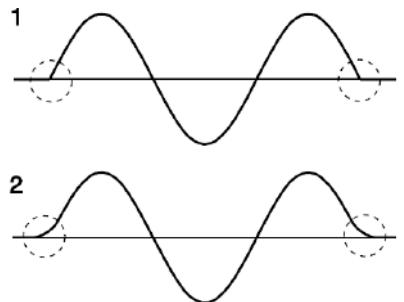
electric field, it must have a more complicated wave pattern. Let's consider the example of an electron being accelerated by the electron gun at the back of a TV tube. Newton's laws are not useful, because they implicitly assume that the path taken by the particle is a meaningful concept. Conservation of energy is still valid in quantum physics, however. In terms of energy, the electron is moving from a region of low voltage into a region of higher voltage. Since its charge is negative, it loses electrical energy by moving to a higher voltage, so its kinetic energy increases. As its electrical energy goes down, its kinetic energy goes up by an equal amount, keeping the total energy constant. Increasing kinetic energy implies a growing momentum, and therefore a shortening wavelength, λ .

The wavefunction as a whole does not have a single well-defined wavelength, but the wave changes so gradually that if you only look at a small part of it you can still pick out a wavelength and relate it to the momentum and energy. (The picture actually exaggerates by many orders of magnitude the rate at which the wavelength changes.)

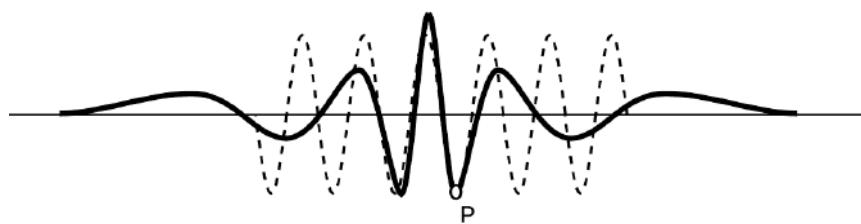
But what if the electric field was stronger? The electric field in an old-fashioned vacuum tube TV screen is only $\sim 10^5 \text{ N/C}$, but the electric field within an atom is more like 10^{12} N/C . In figure 1, the wavelength changes so rapidly that there is nothing that looks like a sine wave at all. We could get a rough idea of the wavelength in a given region by measuring the distance between two peaks, but that would only be a rough approximation. Suppose we want to know the wavelength at point P . The trick is to construct a sine wave, like the one shown with the dashed line, which matches the curvature of the actual wavefunction as closely as possible near P . The sine wave that matches as well as possible is called the "osculating" curve, from a Latin word meaning "to kiss." The wavelength of the osculating curve is the wavelength that will relate correctly to conservation of energy.



j / An electron in a gentle electric field gradually shortens its wavelength as it gains energy. (As discussed on p. 896, it is actually not quite correct to graph the wavefunction of an electron as a real number unless it is a standing wave, which isn't the case here.)



k / The wavefunction's tails go where classical physics says they shouldn't.



l / A typical wavefunction of an electron in an atom (heavy curve) and the osculating sine wave (dashed curve) that matches its curvature at point P .

Tunneling

We implicitly assumed that the particle-in-a-box wavefunction would cut off abruptly at the sides of the box, $k/1$, but that would be unphysical. A kink has infinite curvature, and curvature is related to energy, so it can't be infinite. A physically realistic wavefunction must always “tail off” gradually, $k/2$. In classical physics, a particle can never enter a region in which its interaction energy U would be greater than the amount of energy it has available. But in quantum physics the wavefunction will always have a tail that reaches into the classically forbidden region. If it was not for this effect, called tunneling, the fusion reactions that power the sun would not occur due to the high electrical energy nuclei need in order to get close together! Tunneling is discussed in more detail in the following subsection.

13.3.6 The Schrödinger equation

In subsection 13.3.5 we were able to apply conservation of energy to an electron’s wavefunction, but only by using the clumsy graphical technique of osculating sine waves as a measure of the wave’s curvature. You have learned a more convenient measure of curvature in calculus: the second derivative. To relate the two approaches, we take the second derivative of a sine wave:

$$\begin{aligned}\frac{d^2}{dx^2} \sin(2\pi x/\lambda) &= \frac{d}{dx} \left(\frac{2\pi}{\lambda} \cos \frac{2\pi x}{\lambda} \right) \\ &= - \left(\frac{2\pi}{\lambda} \right)^2 \sin \frac{2\pi x}{\lambda}\end{aligned}$$

Taking the second derivative gives us back the same function, but with a minus sign and a constant out in front that is related to the wavelength. We can thus relate the second derivative to the osculating wavelength:

$$[1] \quad \frac{d^2 \Psi}{dx^2} = - \left(\frac{2\pi}{\lambda} \right)^2 \Psi$$

This could be solved for λ in terms of Ψ , but it will turn out below to be more convenient to leave it in this form.

Applying this to conservation of energy, we have

$$\begin{aligned}[2] \quad E &= K + U \\ &= \frac{p^2}{2m} + U \\ &= \frac{(h/\lambda)^2}{2m} + U\end{aligned}$$

Note that both equation [1] and equation [2] have λ^2 in the denominator. We can simplify our algebra by multiplying both sides of equation [2] by Ψ to make it look more like equation [1]:

$$\begin{aligned} E \cdot \Psi &= \frac{(h/\lambda)^2}{2m} \Psi + U \cdot \Psi \\ &= \frac{1}{2m} \left(\frac{h}{2\pi} \right)^2 \left(\frac{2\pi}{\lambda} \right)^2 \Psi + U \cdot \Psi \\ &= -\frac{1}{2m} \left(\frac{h}{2\pi} \right)^2 \frac{d^2 \Psi}{dx^2} + U \cdot \Psi \end{aligned}$$

Further simplification is achieved by using the symbol \hbar (h with a slash through it, read “h-bar”) as an abbreviation for $h/2\pi$. We then have the important result known as the **Schrödinger equation**:

$$E \cdot \Psi = -\frac{\hbar^2}{2m} \frac{d^2 \Psi}{dx^2} + U \cdot \Psi$$

(Actually this is a simplified version of the Schrödinger equation, applying only to standing waves in one dimension.) Physically it is a statement of conservation of energy. The total energy E must be constant, so the equation tells us that a change in interaction energy U must be accompanied by a change in the curvature of the wavefunction. This change in curvature relates to a change in wavelength, which corresponds to a change in momentum and kinetic energy.

self-check G

Considering the assumptions that were made in deriving the Schrödinger equation, would it be correct to apply it to a photon? To an electron moving at relativistic speeds?

▷ Answer, p. 1067

1067

Usually we know right off the bat how U depends on x , so the basic mathematical problem of quantum physics is to find a function $\Psi(x)$ that satisfies the Schrödinger equation for a given interaction-energy function $U(x)$. An equation, such as the Schrödinger equation, that specifies a relationship between a function and its derivatives is known as a differential equation.

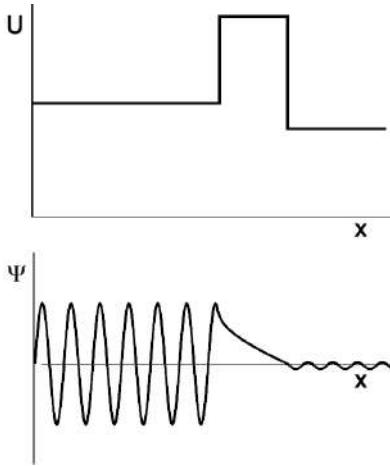
The detailed study of the solution of the Schrödinger equation is beyond the scope of this book, but we can gain some important insights by considering the easiest version of the Schrödinger equation, in which the interaction energy U is constant. We can then rearrange the Schrödinger equation as follows:

$$\frac{d^2 \Psi}{dx^2} = \frac{2m(U - E)}{\hbar^2} \Psi,$$

which boils down to

$$\frac{d^2 \Psi}{dx^2} = a\Psi,$$

where, according to our assumptions, a is independent of x . We need to find a function whose second derivative is the same as the original function except for a multiplicative constant. The only functions with this property are sine waves and exponentials:



m / Tunneling through a barrier. (As discussed on p. 896, it is actually not quite correct to graph the wavefunction of an electron as a real number unless it is a standing wave, which isn't the case here.)

$$\begin{aligned}\frac{d^2}{dx^2} [q \sin(rx + s)] &= -qr^2 \sin(rx + s) \\ \frac{d^2}{dx^2} [qe^{rx+s}] &= qr^2 e^{rx+s}\end{aligned}$$

The sine wave gives negative values of a , $a = -r^2$, and the exponential gives positive ones, $a = r^2$. The former applies to the classically allowed region with $U < E$.

This leads us to a quantitative calculation of the tunneling effect discussed briefly in the preceding subsection. The wavefunction evidently tails off exponentially in the classically forbidden region. Suppose, as shown in figure m, a wave-particle traveling to the right encounters a barrier that it is classically forbidden to enter. Although the form of the Schrödinger equation we're using technically does not apply to traveling waves (because it makes no reference to time), it turns out that we can still use it to make a reasonable calculation of the probability that the particle will make it through the barrier. If we let the barrier's width be w , then the ratio of the wavefunction on the left side of the barrier to the wavefunction on the right is

$$\frac{qe^{rx+s}}{qe^{r(x+w)+s}} = e^{-rw}.$$

Probabilities are proportional to the squares of wavefunctions, so

the probability of making it through the barrier is

$$P = e^{-2rw}$$

$$= \exp\left(-\frac{2w}{\hbar}\sqrt{2m(U-E)}\right).$$

self-check H

If we were to apply this equation to find the probability that a person can walk through a wall, what would the small value of Planck's constant imply?

▷ Answer, p. 1067

Tunneling in alpha decay

example 17

Naively, we would expect alpha decay to be a very fast process. The typical speeds of neutrons and protons inside a nucleus are extremely high (see problem 20). If we imagine an alpha particle coalescing out of neutrons and protons inside the nucleus, then at the typical speeds we're talking about, it takes a ridiculously small amount of time for them to reach the surface and try to escape. Clattering back and forth inside the nucleus, we could imagine them making a vast number of these "escape attempts" every second.

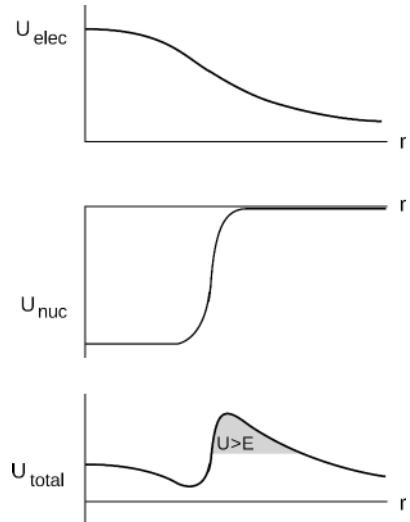
Consider figure n, however, which shows the interaction energy for an alpha particle escaping from a nucleus. The electrical energy is $kq_1 q_2 / r$ when the alpha is outside the nucleus, while its variation inside the nucleus has the shape of a parabola, as a consequence of the shell theorem. The nuclear energy is constant when the alpha is inside the nucleus, because the forces from all the neighboring neutrons and protons cancel out; it rises sharply near the surface, and flattens out to zero over a distance of ~ 1 fm, which is the maximum distance scale at which the strong force can operate. There is a classically forbidden region immediately outside the nucleus, so the alpha particle can only escape by quantum mechanical tunneling. (It's true, but somewhat counterintuitive, that a *repulsive* electrical force can make it more difficult for the alpha to get *out*.)

In reality, alpha-decay half-lives are often extremely long — sometimes billions of years — because the tunneling probability is so small. Although the shape of the barrier is not a rectangle, the equation for the tunneling probability on page 909 can still be used as a rough guide to our thinking. Essentially the tunneling probability is so small because $U - E$ is fairly big, typically about 30 MeV at the peak of the barrier.

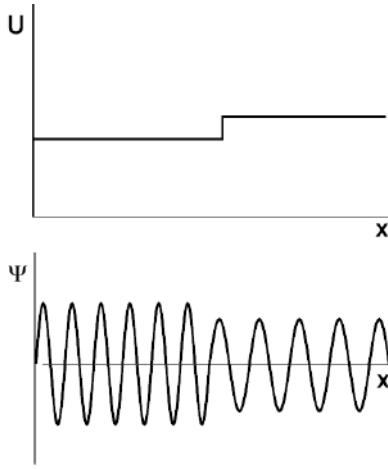
The correspondence principle for $E > U$

example 18

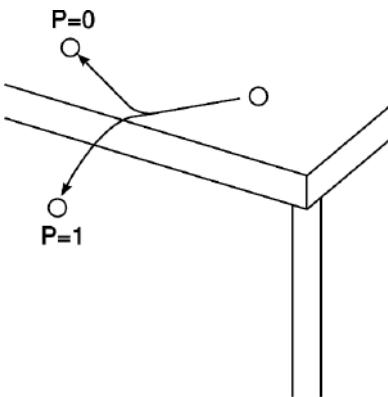
The correspondence principle demands that in the classical limit $\hbar \rightarrow 0$, we recover the correct result for a particle encountering a barrier U , for both $E < U$ and $E > U$. The $E < U$ case was analyzed in self-check H on p. 909. In the remainder of this example, we analyze $E > U$, which turns out to be a little trickier.



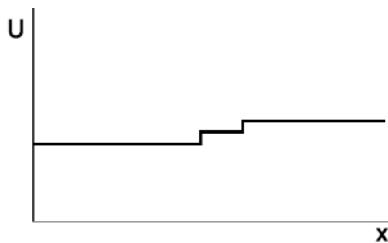
n / The electrical, nuclear, and total interaction energies for an alpha particle escaping from a nucleus.



o / A particle encounters a step of height $U < E$ in the interaction energy. Both sides are classically allowed. A reflected wave exists, but is not shown in the figure.



p / The marble has zero probability of being reflected from the edge of the table. (This example has $U < 0$, not $U > 0$ as in figures o and q).



q / Making the step more gradual reduces the probability of reflection.

The particle has enough energy to get over the barrier, and the classical result is that it continues forward at a different speed (a reduced speed if $U > 0$, or an increased one if $U < 0$), then regains its original speed as it emerges from the other side. What happens quantum-mechanically in this case? We would like to get a “tunneling” probability of 1 in the classical limit. The expression derived on p. 909, however, doesn’t apply here, because it was derived under the assumption that the wavefunction inside the barrier was an exponential; in the classically allowed case, the barrier isn’t classically forbidden, and the wavefunction inside it is a sine wave.

We can simplify things a little by letting the width w of the barrier go to infinity. Classically, after all, there is no possibility that the particle will turn around, no matter how wide the barrier. We then have the situation shown in figure o.⁶

The analysis is similar to that for any other wave being partially reflected at the boundary between two regions where its velocity differs, and the result is the same as the one found on p. 381. (There are some technical differences, which don’t turn out to matter. This is discussed in more detail on p. 974.) The ratio of the amplitude of the reflected wave to that of the incident wave is $R = (v_2 - v_1)/(v_2 + v_1)$. The probability of reflection is R^2 . (Counterintuitively, R^2 is nonzero even if $U < 0$, i.e., $v_2 > v_1$.)

This seems to violate the correspondence principle. There is no m or h anywhere in the result, so we seem to have the result that, even classically, the marble in figure p can be reflected!

The solution to this paradox is that the step in figure o was taken to be completely abrupt — an idealized mathematical discontinuity. Suppose we make the transition a little more gradual, as in figure q. As shown in problem 17 on p. 395, this reduces the amplitude with which a wave is reflected. By smoothing out the step more and more, we continue to reduce the probability of reflection, until finally we arrive at a barrier shaped like a smooth ramp. More detailed calculations show that this results in zero reflection in the limit where the width of the ramp is large compared to the wavelength.

Beta decay: a push or pull on the way out the door example 19
The nucleus ^{64}Cu undergoes β^+ and β^- decay with similar probabilities and energies. Each of these decays releases a fixed amount of energy Q due to the difference in mass between the parent nucleus and the decay products. This energy is shared randomly between the beta and the neutrino. In experiments, the beta’s energy is easily measured, while the neutrino flies off without interacting. Figure r shows the energy spectrum of the β^+ and

⁶As in several previous examples, we cheat by representing a traveling wave as a real-valued function. See pp. 896 and 974.

β^- in these decays.⁷ There is a relatively high probability for the beta and neutrino each to carry off roughly half the kinetic energy, the reason being identical to the kind of phase-space argument discussed in sec. 5.4.2, p. 328. Therefore in each case we get a bell-shaped curve stretching from 0 up to the energy Q , with Q being slightly different in the two cases.

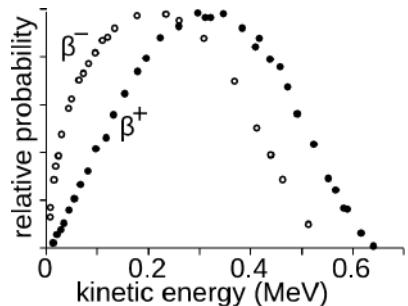
So we expect the two bell curves to look almost the same except for a slight rescaling of the horizontal axis. Yes — but we also see markedly different behavior at low energies. At very low energies, there is almost no chance to see a β^+ with very low energy, but quite a high probability for a β^- .

We could try to explain this difference in terms of the release of electrical energy. The β^+ is repelled by the nucleus, so it gets an extra push on the way out the door. A β^- should be held back as it exits, and so should lose some energy. The bell curves should be shifted up and down in energy relative to one another, as observed.

But if we try to estimate this energy shift using classical physics, we come out with a wildly incorrect answer. This would be a process in which the beta and neutrino are released in a point-like event inside the nucleus. The radius r of the ^{64}Cu nucleus is on the order of 4 fm ($1 \text{ fm} = 10^{-15} \text{ m}$). Therefore the energy lost or gained by the β^+ or β^- on the way out would be $U \sim kZ\epsilon^2/r \sim 10 \text{ MeV}$. The actual shift is much smaller.

To understand what's really going on, we need quantum mechanics. A beta in the observed energy range has a wavelength of about 2000 fm, which is hundreds of times greater than the size of the nucleus. Therefore the beta cannot be much better localized than that when it is emitted. This means that we should really use something more like $r \sim 500 \text{ fm}$ (a quarter of a wavelength) in our calculation of the electrical energy. This gives $U \sim 0.08 \text{ MeV}$, which is about the right order of magnitude compared to observation.

A byproduct of this analysis is that a β^+ is always emitted within the classically forbidden region, and then has to tunnel out through the barrier. As in example 17, we have the counterintuitive fact about quantum mechanics that a repulsive force can *hinder* the escape of a particle.

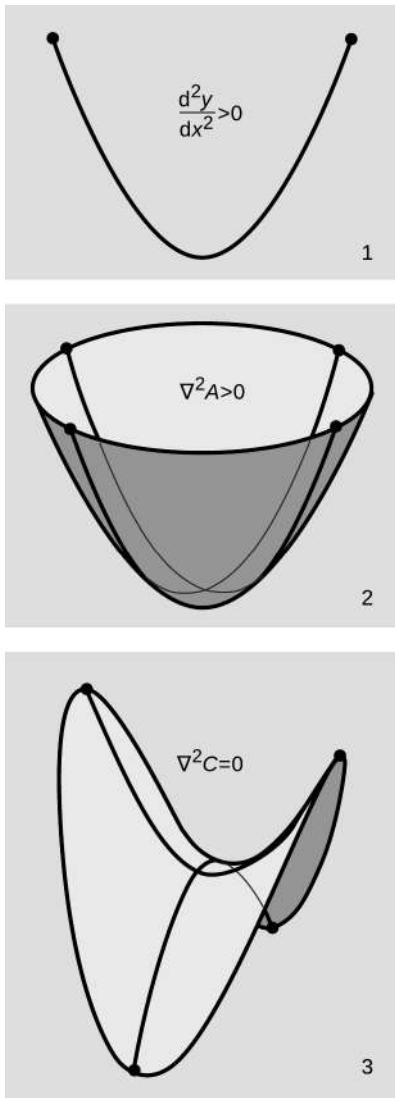


r / β^+ and β^- spectra of ^{64}Cu .

⁷Redrawn from Cook and Langer, 1948.

Three dimensions

For simplicity, we've been considering the Schrödinger equation in one dimension, so that Ψ is only a function of x , and has units of $m^{-1/2}$ rather than $m^{-3/2}$. Since the Schrödinger equation is a statement of conservation of energy, and energy is a scalar, the generalization to three dimensions isn't particularly complicated. The total energy term $E \cdot \Psi$ and the interaction energy term $U \cdot \Psi$ involve nothing but scalars, and don't need to be changed at all. In the kinetic energy term, however, we're essentially basing our computation of the kinetic energy on the squared magnitude of the momentum, p_x^2 , and in three dimensions this would clearly have to be generalized to $p_x^2 + p_y^2 + p_z^2$. The obvious way to achieve this is to replace the second derivative $d^2\Psi/dx^2$ with the sum $\partial^2\Psi/\partial x^2 + \partial^2\Psi/\partial y^2 + \partial^2\Psi/\partial z^2$. Here the partial derivative symbol ∂ , introduced on page 220, indicates that when differentiating with respect to a particular variable, the other variables are to be considered as constants. This operation on the function Ψ is notated $\nabla^2\Psi$, and the derivative-like operator $\nabla^2 = \partial^2/\partial x^2 + \partial^2/\partial y^2 + \partial^2/\partial z^2$ is called the Laplacian, and was introduced briefly on p. 657. It occurs elsewhere in physics. For example, in classical electrostatics, the voltage in a region of vacuum must be a solution of the equation $\nabla^2V = 0$. Like the second derivative, the Laplacian is essentially a measure of curvature. Or, as shown in figure s, we can think of it as a measure of how much the value of a function at a certain point differs from the average of its value on nearby points.



s / 1. The one-dimensional version of the Laplacian is the second derivative. It is positive here because the average of the two nearby points is greater than the value at the center. 2. The Laplacian of the function A in example 20 is positive because the average of the four nearby points along the perpendicular axes is greater than the function's value at the center. 3. $\nabla^2C = 0$. The average is the same as the value at the center.

Examples of the Laplacian in two dimensions *example 20*

- ▷ Compute the Laplacians of the following functions in two dimensions, and interpret them: $A = x^2 + y^2$, $B = -x^2 - y^2$, $C = x^2 - y^2$.
- ▷ The first derivative of function A with respect to x is $\partial A/\partial x = 2x$. Since y is treated as a constant in the computation of the partial derivative $\partial/\partial x$, the second term goes away. The second derivative of A with respect to x is $\partial^2 A/\partial x^2 = 2$. Similarly we have $\partial^2 A/\partial y^2 = 2$, so $\nabla^2 A = 4$.

All derivative operators, including ∇^2 , have the linear property that multiplying the input function by a constant just multiplies the output function by the same constant. Since $B = -A$, and we have $\nabla^2 B = -4$.

For function C , the x term contributes a second derivative of 2, but the y term contributes -2 , so $\nabla^2 C = 0$.

The interpretation of the positive sign in $\nabla^2 A = 4$ is that A 's graph is shaped like a trophy cup, and the cup is concave up. $\nabla^2 B < 0$ is because B is concave down. Function C is shaped like a saddle. Since its curvature along one axis is concave up, but the curvature along the other is down and equal in magnitude, the

function is considered to have zero concavity over all.

A classically allowed region with constant U *example 21*

In a classically allowed region with constant U , we expect the solutions to the Schrödinger equation to be sine waves. A sine wave in three dimensions has the form

$$\Psi = \sin(k_x x + k_y y + k_z z).$$

When we compute $\partial^2\Psi/\partial x^2$, double differentiation of sin gives $-\sin$, and the chain rule brings out a factor of k_x^2 . Applying all three second derivative operators, we get

$$\begin{aligned}\nabla^2\Psi &= (-k_x^2 - k_y^2 - k_z^2) \sin(k_x x + k_y y + k_z z) \\ &= -(k_x^2 + k_y^2 + k_z^2) \Psi.\end{aligned}$$

The Schrödinger equation gives

$$\begin{aligned}E \cdot \Psi &= -\frac{\hbar^2}{2m} \nabla^2\Psi + U \cdot \Psi \\ &= -\frac{\hbar^2}{2m} \cdot -(k_x^2 + k_y^2 + k_z^2) \Psi + U \cdot \Psi \\ E - U &= \frac{\hbar^2}{2m} (k_x^2 + k_y^2 + k_z^2),\end{aligned}$$

which can be satisfied since we're in a classically allowed region with $E - U > 0$, and the right-hand side is manifestly positive.

Use of complex numbers

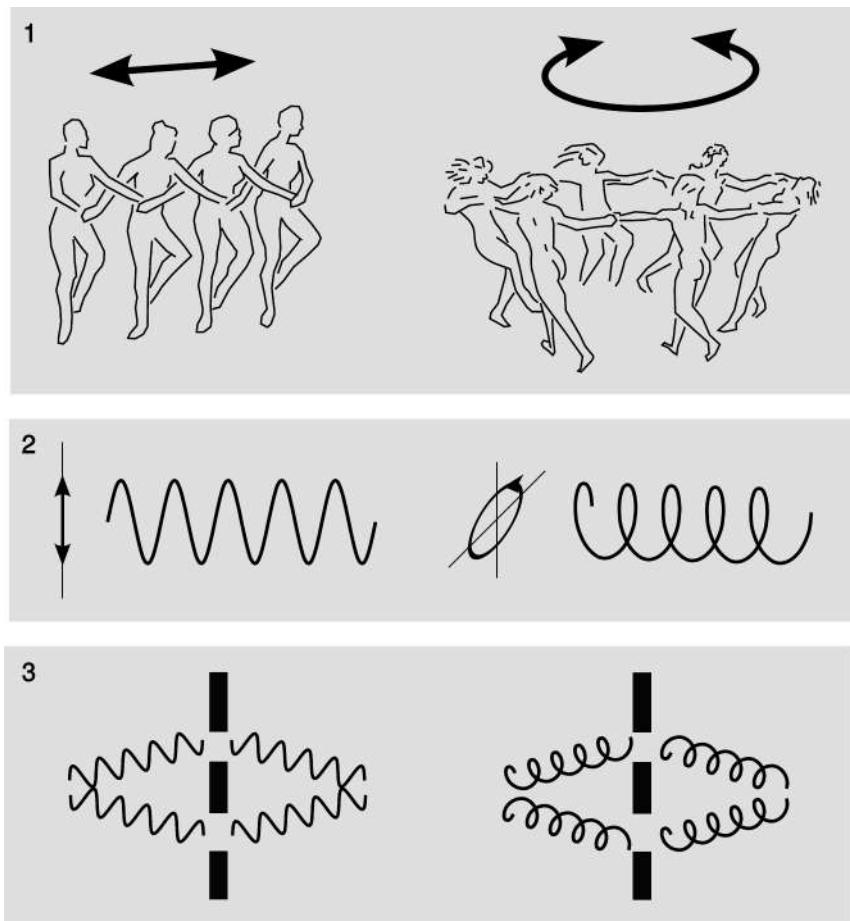
In a classically forbidden region, a particle's total energy, $U + K$, is less than its U , so its K must be negative. If we want to keep believing in the equation $K = p^2/2m$, then apparently the momentum of the particle is the square root of a negative number. This is a symptom of the fact that the Schrödinger equation fails to describe all of nature unless the wavefunction and various other quantities are allowed to be complex numbers. In particular it is not possible to describe traveling waves correctly without using complex wavefunctions. Complex numbers were reviewed in subsection 10.5.5, p. 627.

This may seem like nonsense, since real numbers are the only ones that are, well, real! Quantum mechanics can always be related to the real world, however, because its structure is such that the results of measurements always come out to be real numbers. For example, we may describe an electron as having non-real momentum in classically forbidden regions, but its average momentum will always come out to be real (the imaginary parts average out to zero), and it can never transfer a non-real quantity of momentum to another particle.

t / 1. Oscillations can go back and forth, but it's also possible for them to move along a path that bites its own tail, like a circle. Photons act like one, electrons like the other.

2. Back-and-forth oscillations can naturally be described by a segment taken from the real number line, and we visualize the corresponding type of wave as a sine wave. Oscillations around a closed path relate more naturally to the complex number system. The complex number system has rotation built into its structure, e.g., the sequence $1, i, i^2, i^3, \dots$ rotates around the unit circle in 90-degree increments.

3. The double slit experiment embodies the one and only mystery of quantum physics. Either type of wave can undergo double-slit interference.



A complete investigation of these issues is beyond the scope of this book, and this is why we have normally limited ourselves to standing waves, which can be described with real-valued wavefunctions. Figure t gives a visual depiction of the difference between real and complex wavefunctions. The following remarks may also be helpful.

Neither of the graphs in t/2 should be interpreted as a path traveled by something. This isn't anything mystical about quantum physics. It's just an ordinary fact about waves, which we first encountered in subsection 6.1.1, p. 354, where we saw the distinction between the motion of a wave and the motion of a wave pattern. In both examples in t/2, the wave pattern is moving in a straight line to the right.

The helical graph in t/2 shows a complex wavefunction whose value rotates around a circle in the complex plane with a frequency f related to its energy by $E = hf$. As it does so, its squared magnitude $|\Psi|^2$ stays the same, so the corresponding probability stays constant. Which direction does it rotate? This direction is purely a matter of convention, since the distinction between the symbols i and $-i$ is arbitrary — both are equally valid as square roots of -1 . We can,

for example, arbitrarily say that electrons with positive energies have wavefunctions whose phases rotate counterclockwise, and as long as we follow that rule consistently within a given calculation, everything will work. Note that it is not possible to define anything like a right-hand rule here, because the complex plane shown in the right-hand side of t/2 doesn't represent two dimensions of physical space; unlike a screw going into a piece of wood, an electron doesn't have a direction of rotation that depends on its direction of travel.

Superposition of complex wavefunctions *example 22*

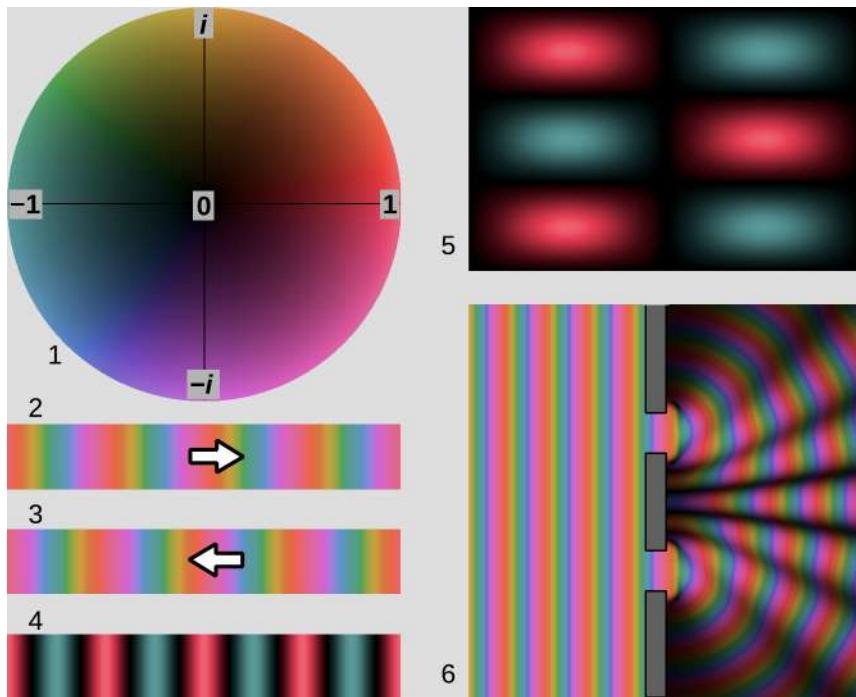
- ▷ The right side of figure t/3 is a cartoonish representation of double-slit interference; it depicts the situation at the center, where symmetry guarantees that the interference is constructive. Suppose that at some off-center point, the two wavefunctions being superposed are $\Psi_1 = b$ and $\Psi_2 = bi$, where b is a real number with units. Compare the probability of finding the electron at this position with what it would have been if the superposition had been purely constructive, $b + b = 2b$.
- ▷ The probability per unit volume is proportional to the square of the magnitude of the total wavefunction, so we have

$$\frac{P_{\text{off center}}}{P_{\text{center}}} = \frac{|b + bi|^2}{|b + b|^2} = \frac{1^2 + 1^2}{2^2 + 0^2} = \frac{1}{2}.$$

Figure u shows a method for visualizing complex wavefunctions. The idea is to use colors to represent complex numbers, according to the arbitrary convention defined in figure u/1. Brightness indicates magnitude, and the rainbow hue shows the argument. Because this representation can't be understood in a black and white printed book, the figure is also reproduced on the back cover of printed copies. To avoid any confusion, note that the use of rainbow colors does not mean that we are representing actual visible light. In fact, we will be using these visual conventions to represent the wavefunctions of a material particle such as an electron. It is arbitrary that we use red for positive real numbers and blue-green for negative numbers, and that we pick a handedness for the diagram such that going from red toward yellow means going counterclockwise. Although physically the rainbow is a linear spectrum, we are not representing physical colors here, and we are exploiting the fact that the human brain tends to perceive color as a circle rather than a line, with violet and red being perceptually similar. One of the limitations of this representation is that brightness is limited, so we can't represent complex numbers with arbitrarily large magnitudes.

Figure u/2 shows a traveling wave as it propagates to the right. The standard convention in physics is that for a wave moving in a certain direction, the phase in the forward direction is farther counterclockwise in the complex plane, and you can verify for yourself

u / 1. A representation of complex numbers using color and brightness. 2. A wave traveling toward the right. 3. A wave traveling toward the left. 4. A standing wave formed by superposition of waves 2 and 3. 5. A two-dimensional standing wave. 6. A double-slit diffraction pattern.



that this is the case by comparing with the convention defined by u/1. The function being plotted here is $\Psi = e^{ikx}$, where $k = 2\pi/\lambda$ is the spatial analog of frequency, with an extra factor of 2π for convenience. For the use of the complex exponential, see sec. 10.5.6, p .629; it simply represents a point on the unit circle in the complex plane. The wavelength λ is a constant and can be measured, for example, from one yellow point to the next. The wavelength is *not* different at different points on the figure, because we are using the colors merely as a visual encoding of the complex numbers — so, for example, a red point on the figure is not a point where the wave has a longer wavelength than it does at a blue point.

Figure u/3 represents a wave traveling to the left.

Figure u/4 shows a standing wave created by superimposing the traveling waves from u/2 and u/3, $\Psi_4 = (\Psi_2 + \Psi_3)/2$. (The reason for the factor of 2 is simply that otherwise some portions of Ψ_4 would have magnitudes too great to be represented using the available range of brightness.) All points on this wave have real values, represented by red and blue-green. We made the superposition real by an appropriate choice of the phases of Ψ_2 and Ψ_3 . This is always possible to do when we have a standing wave, but it is *only* possible for a standing wave, and this is the reason for all of the disclaimers in the captions of previous figures in which I took the liberty of representing a traveling wave as a sine-wave graph.

Figure u/5 shows a two-dimensional standing wave of a particle in a box, and u/6 shows a double-slit interference pattern. (In the latter, I've cheated by making the amplitude of the wave on the

right-hand half of the picture much greater than it would actually be.)

A paradox resolved

example 23

Consider the following paradox. Suppose we have an electron that is traveling wave, and its wavefunction looks like a wave-train consisting of 5 cycles of a sine wave. Call the distance between the leading and trailing edges of the wave-train L , so that $\lambda = L/5$. By sketching the wave, you can easily check that there are 11 points where its value equals zero. Therefore at a particular moment in time, there are 11 points where a detector has zero probability of detecting the electron.

But now consider how this would look in a frame of reference where the electron is moving more slowly, at one fifth of the speed we saw in the original frame. In this frame, L is the same, but λ is five times greater, because $\lambda = h/p$. Therefore in this frame we see only one cycle in the wave-train. Now there are only 3 points where the probability of detection is zero. But how can this be? All observers, regardless of their frames of reference, should agree on whether a particular detector detects the electron.

The resolution to this paradox is that it starts from the assumption that we can depict a traveling wave as a real-valued sine wave, which is zero in certain places. Actually, we can't. It has to be a complex number with a rotating phase angle in the complex plane, as in figure u/2, and a *constant* magnitude.

Linearity of the Schrödinger equation

Some mathematical relationships and operations are *linear*, and some are not. For example, $2 \times (3+2)$ is the same as $2 \times 3 + 2 \times 2$, but $\sqrt{1+1} \neq \sqrt{1} + \sqrt{1}$. Differentiation is a linear operation, $(f+g)' = f' + g'$. The Schrödinger equation is built out of derivatives, so it is linear as well. That is, if Ψ_1 and Ψ_2 are both solutions of the Schrödinger equation, then so is $\Psi_1 + \Psi_2$. Linearity normally implies linearity with respect both to addition and to multiplication by a scalar. For example, if Ψ is a solution, then so is $\Psi + \Psi + \Psi$, which is the same as 3Ψ .

Linearity guarantees that the phase of a wavefunction makes no difference as to its validity as a solution to the Schrödinger equation. If $\sin kx$ is a solution, then so is the sine wave $-\sin kx$ with the opposite phase. This fact is logically interdependent with the fact that, as discussed on p. 896, the phase of a wavefunction is unobservable. For measuring devices and humans are material objects that can be described by wavefunctions. So suppose, for example, that we flip the phase of all the particles inside the entire laboratory. By linearity, the evolution of this measurement process is still a valid solution of the Schrödinger equation.

The Schrödinger equation is a wave equation, and its linearity

implies that the waves obey the principle of superposition. In most cases in nature, we find that the principle of superposition for waves is at best an approximation. For example, if the amplitude of a tsunami is so huge that the trough of the wave reaches all the way down to the ocean floor, exposing the rocks and sand as it passes overhead, then clearly there is no way to double the amplitude of the wave and still get something that obeys the laws of physics. Even at less extreme amplitudes, superposition is only an approximation for water waves, and so for example it is only approximately true that when two sets of ripples intersect on the surface of a pond, they pass through without “seeing” each other.

It is therefore natural to ask whether the apparent linearity of the Schrödinger equation is only an approximation to some more precise, nonlinear theory. This is not currently believed to be the case. If we are to make sense of Schrödinger’s cat (sec. 13.2.4, p. 887), then the experimenter who sees a live cat and the one who sees a dead cat must remain oblivious to their other selves, like the ripples on the pond that intersect without “seeing” each other. Attempts to create slightly nonlinear versions of standard quantum mechanics have been shown to have implausible physical properties, such as allowing the propagation of signals faster than c . (This is known as Gisin’s theorem. The original paper, “Weinberg’s non-linear quantum mechanics and supraluminal communications,” is surprisingly readable and nonmathematical.)

If you have had a course in linear algebra, then it is worth noting that the linearity of the Schrödinger equation allows us to talk about its solutions as vectors in a vector space. For example, if Ψ_1 represents an unstable nucleus that has not yet gamma decayed, and Ψ_2 is its state after the decay, then any superposition $\alpha\Psi_1 + \beta\Psi_2$, with real or complex coefficients α and β , is a possible wavefunction, and we can notate this as a vector, $\langle\alpha, \beta\rangle$, in a two-dimensional vector space.

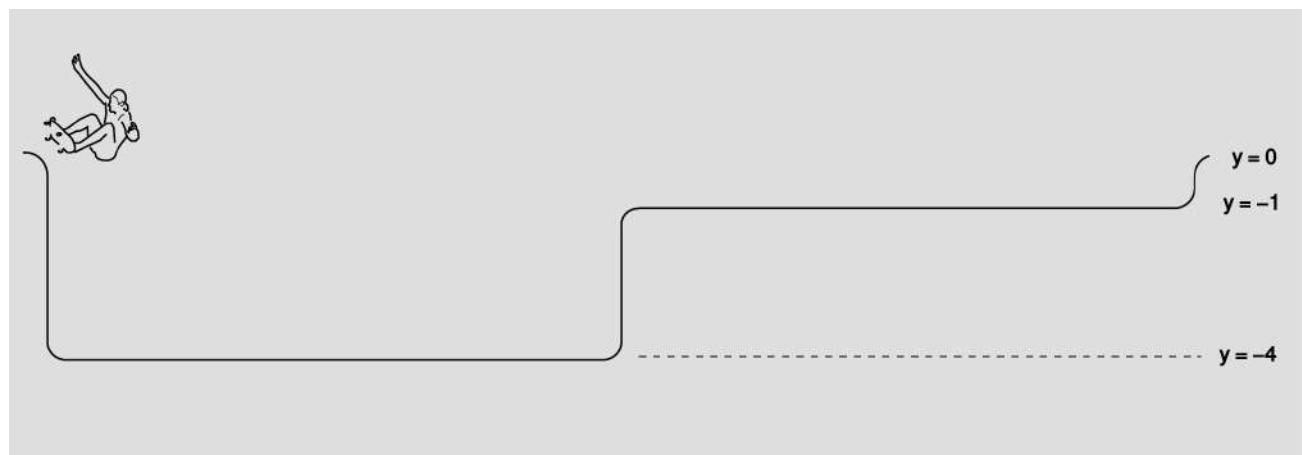
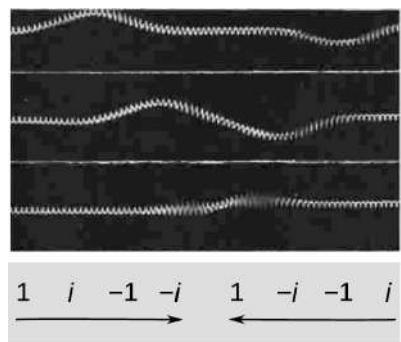
Discussion Questions

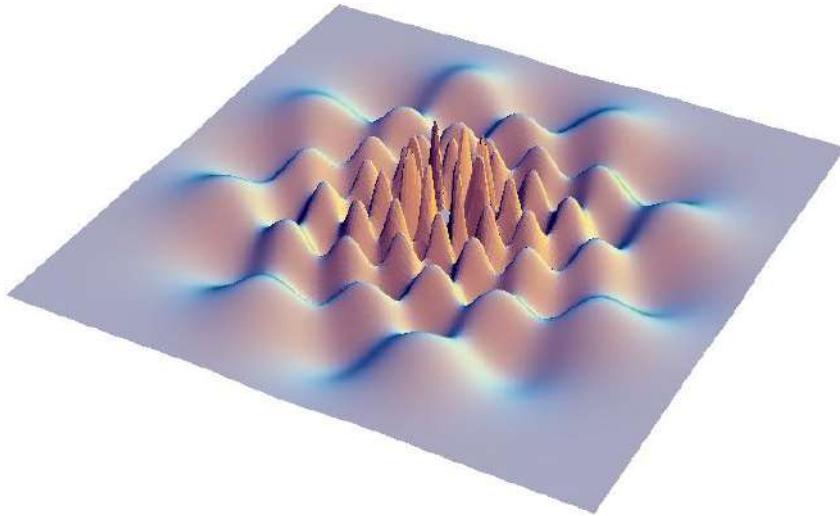
A The zero level of interaction energy U is arbitrary, e.g., it’s equally valid to pick the zero of gravitational energy to be on the floor of your lab or at the ceiling. Suppose we’re doing the double-slit experiment, t/3, with electrons. We define the zero-level of U so that the total energy $E = U + K$ of each electron is positive. and we observe a certain interference pattern like the one in figure i on p. 880. What happens if we then redefine the zero-level of U so that the electrons have $E < 0$?

B The top panel of the figure shows a series of snapshots in the motion of two pulses on a coil spring, one negative and one positive, as they move toward one another and superpose. The final image is very close to the moment at which the two pulses cancel completely. The following discussion is simpler if we consider infinite sine waves rather than pulses. How can the cancellation of two such mechanical waves be reconciled with conservation of energy? What about the case of colliding electromagnetic waves?

Quantum-mechanically, the issue isn't conservation of energy, it's conservation of probability, i.e., if there's initially a 100% probability that a particle exists somewhere, we don't want the probability to be more than or less than 100% at some later time. What happens when the colliding waves have real-valued wavefunctions Ψ ? Now consider the sketches of complex-valued wave pulses shown in the bottom panel of the figure as they are getting ready to collide.

C The figure shows a skateboarder tipping over into a swimming pool with zero initial kinetic energy. There is no friction, the corners are smooth enough to allow the skater to pass over the smoothly, and the vertical distances are small enough so that negligible time is required for the vertical parts of the motion. The pool is divided into a deep end and a shallow end. Their widths are equal. The deep end is four times deeper. (1) Classically, compare the skater's velocity in the left and right regions, and infer the probability of finding the skater in either of the two halves if an observer peeks at a random moment. (2) Quantum-mechanically, this could be a one-dimensional model of an electron shared between two atoms in a diatomic molecule. Compare the electron's kinetic energies, momenta, and wavelengths in the two sides. For simplicity, let's assume that there is no tunneling into the classically forbidden regions. What is the simplest standing-wave pattern that you can draw, and what are the probabilities of finding the electron in one side or the other? Does this obey the correspondence principle?





13.4 The atom

You can learn a lot by taking a car engine apart, but you will have learned a lot more if you can put it all back together again and make it run. Half the job of reductionism is to break nature down into its smallest parts and understand the rules those parts obey. The second half is to show how those parts go together, and that is our goal in this chapter. We have seen how certain features of all atoms can be explained on a generic basis in terms of the properties of bound states, but this kind of argument clearly cannot tell us any details of the behavior of an atom or explain why one atom acts differently from another.

The biggest embarrassment for reductionists is that the job of putting things back together job is usually much harder than the taking them apart. Seventy years after the fundamentals of atomic physics were solved, it is only beginning to be possible to calculate accurately the properties of atoms that have many electrons. Systems consisting of many atoms are even harder. Supercomputer manufacturers point to the folding of large protein molecules as a process whose calculation is just barely feasible with their fastest machines. The goal of this chapter is to give a gentle and visually oriented guide to some of the simpler results about atoms.

13.4.1 Classifying states

We'll focus our attention first on the simplest atom, hydrogen, with one proton and one electron. We know in advance a little of what we should expect for the structure of this atom. Since the electron is bound to the proton by electrical forces, it should display a set of discrete energy states, each corresponding to a certain standing wave pattern. We need to understand what states there are and what their properties are.

What properties should we use to classify the states? The most sensible approach is to use conserved quantities. Energy is one conserved quantity, and we already know to expect each state to have a specific energy. It turns out, however, that energy alone is not sufficient. Different standing wave patterns of the atom can have the same energy.

Momentum is also a conserved quantity, but it is not particularly appropriate for classifying the states of the electron in a hydrogen atom. The reason is that the force between the electron and the proton results in the continual exchange of momentum between them. (Why wasn't this a problem for energy as well? Kinetic energy and momentum are related by $K = p^2/2m$, so the much more massive proton never has very much kinetic energy. We are making an approximation by assuming all the kinetic energy is in the electron, but it is quite a good approximation.)

Angular momentum does help with classification. There is no transfer of angular momentum between the proton and the electron, since the force between them is a center-to-center force, producing no torque.

Like energy, angular momentum is quantized in quantum physics. As an example, consider a quantum wave-particle confined to a circle, like a wave in a circular moat surrounding a castle. A sine wave in such a "quantum moat" cannot have any old wavelength, because an integer number of wavelengths must fit around the circumference, C , of the moat. The larger this integer is, the shorter the wavelength, and a shorter wavelength relates to greater momentum and angular momentum. Since this integer is related to angular momentum, we use the symbol ℓ for it:

$$\lambda = C/\ell$$

The angular momentum is

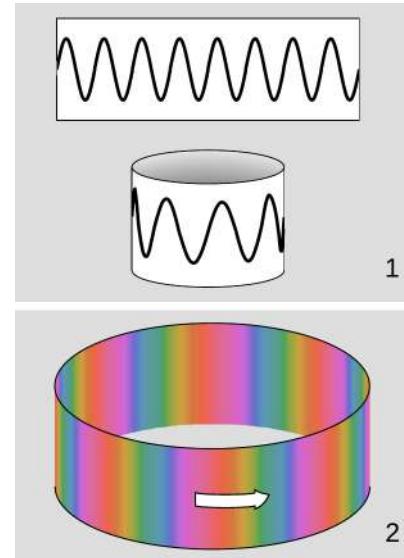
$$L = rp.$$

Here, $r = C/2\pi$, and $p = h/\lambda = h\ell/C$, so

$$\begin{aligned} L &= \frac{C}{2\pi} \cdot \frac{h\ell}{C} \\ &= \frac{h}{2\pi}\ell \end{aligned}$$

In the example of the quantum moat, angular momentum is quantized in units of $h/2\pi$. This makes $h/2\pi$ a pretty important number, so we define the abbreviation $\hbar = h/2\pi$. This symbol is read "h-bar."

In fact, this is a completely general fact in quantum physics, not just a fact about the quantum moat:



a / 1. Eight wavelengths fit around this circle ($\ell = 8$). This is a standing wave. 2. A traveling wave with $\ell = 8$, depicted according to the color conventions defined in figure u, p. 916.

Quantization of angular momentum

The angular momentum of a particle due to its motion through space is quantized in units of \hbar .

self-check 1

What is the angular momentum of the wavefunction shown at the beginning of the section?

▷ Answer, p.

1067

Degeneracy

Comparing the oversimplified figure a/1 with the more realistic depiction in a/2 using complex numbers, we see that there is a direction of rotation, which was drawn as counterclockwise in the figure. As in figures u/2 and u/3 on p. 916, we could have drawn the clockwise version by putting the rainbow colors in the opposite order, i.e., by letting the phase spin in the opposite direction in the complex plane. This feature was hidden in a/1, where in order to get a depiction using real numbers, we had to use a standing wave. A standing wave, however, can be constructed as a superposition of two traveling waves, so the issue was still there, just hidden. We really have *two* quantum-mechanical states here, regardless of whether we use standing waves or traveling waves. If we use standing waves, they are of the form $\sin 8\theta$ and $\cos 8\theta$, while in terms of the traveling waves we have $e^{8i\theta}$ and $e^{-8i\theta}$. By Euler's formula (sec. 10.5.6, p. 629), either traveling wave can be expressed as a superposition of the two standing waves, and vice versa. (Physically, there are not four different states here but two. The situation is a bit like choosing a Cartesian coordinate system in the plane, where we could choose one coordinate system (x, y) , or some other coordinates system (x', y') rotated with respect to the first one; but this does not mean there are four coordinates needed to describe a plane.)

These two states are simplified models of states in an atom, so it's worth thinking about how we could tell, for a real atom, whether the electron had angular momentum in one direction or the other. One technique would be to look at absorption or emission spectra of thin gases, as in example 13 on p. 900. But this only distinguishes states according to their energies, and since these two states have the same kinetic energy, that would not necessarily help. In quantum mechanics, when we have more than one state with the same energy, they are said to be *degenerate*. In our example, the degeneracy of the $\ell = 8$ state is 2. This degeneracy arises from the symmetry of space, which does not distinguish one direction from another. Degeneracies often, but not always, arise from symmetries. (Cf. p. 929.)

If we wanted to distinguish these two degenerate states observationally, one way to do it would be to "lift" the degeneracy by

applying an external magnetic field. Since an electron has an electric charge, it acts like a current loop, and the two states behave like oppositely oriented magnetic dipoles with an additional potential energy $-\mathbf{m} \cdot \mathbf{B}$, which lowers the energy of one state and raises the energy of the other. The existence of the magnetic field breaks the symmetry, which was the reason for the degeneracy.

13.4.2 Three dimensions

Our discussion of quantum-mechanical angular momentum has so far been limited to rotation in a plane, for which we can simply use positive and negative signs to indicate clockwise and counterclockwise directions of rotation. A hydrogen atom, however, is unavoidably three-dimensional. The classical treatment of angular momentum in three-dimensions has been presented in section 4.3; in general, the angular momentum of a particle is defined as the vector cross product $\mathbf{r} \times \mathbf{p}$.

There is a basic problem here: the angular momentum of the electron in a hydrogen atom depends on both its distance \mathbf{r} from the proton and its momentum \mathbf{p} , so in order to know its angular momentum precisely it would seem we would need to know both its position and its momentum simultaneously with good accuracy. This, however, seems forbidden by the Heisenberg uncertainty principle.

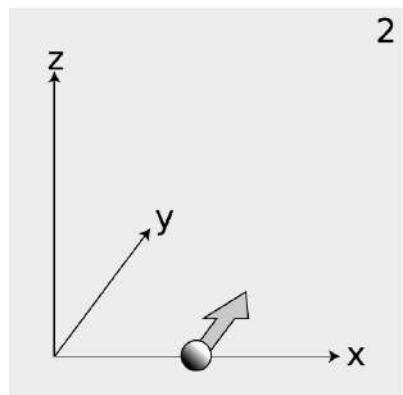
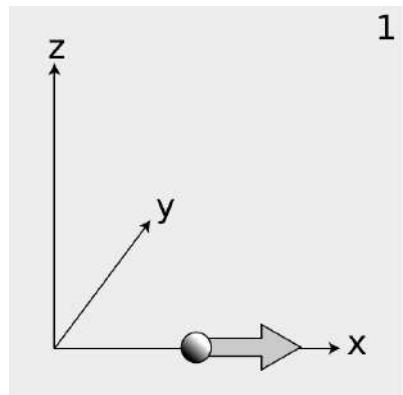
Actually the uncertainty principle does place limits on what can be known about a particle's angular momentum vector, but it does not prevent us from knowing its magnitude as an exact integer multiple of \hbar . The reason is that in three dimensions, there are really three separate uncertainty principles:

$$\begin{aligned}\Delta p_x \Delta x &\gtrsim \hbar \\ \Delta p_y \Delta y &\gtrsim \hbar \\ \Delta p_z \Delta z &\gtrsim \hbar\end{aligned}$$

Now consider a particle, b/1, that is moving along the x axis at position x and with momentum p_x . We may not be able to know both x and p_x with unlimited accuracy, but we can still know the particle's angular momentum about the origin exactly: it is zero, because the particle is moving directly away from the origin.

Suppose, on the other hand, a particle finds itself, b/2, at a position x along the x axis, and it is moving parallel to the y axis with momentum p_y . It has angular momentum xp_y about the z axis, and again we can know its angular momentum with unlimited accuracy, because the uncertainty principle only relates x to p_x and y to p_y . It does not relate x to p_y .

As shown by these examples, the uncertainty principle does not restrict the accuracy of our knowledge of angular momenta as severely as might be imagined. However, it does prevent us from

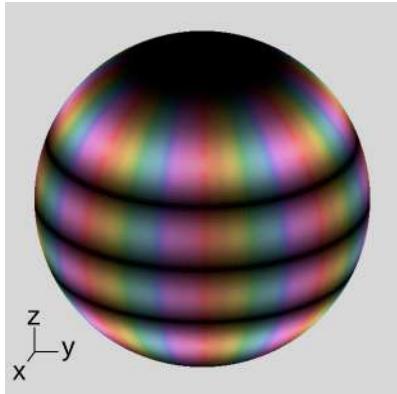


b / Reconciling the uncertainty principle with the definition of angular momentum.

knowing all three components of an angular momentum vector simultaneously. The most general statement about this is the following theorem:

The angular momentum vector in quantum physics

The most that can be known about a (nonzero) orbital angular momentum vector is its magnitude and one of its three vector components. Both are quantized in units of \hbar .



c / A wavefunction on the sphere with $|\mathbf{L}| = 11\hbar$ and $L_z = 8\hbar$, shown using the color conventions defined in figure u, p. 916.

To see why this is true, consider the example wavefunction shown in figure c. This is like the quantum moat of figure a, p. 921, but extended to one more dimension. If we slice the sphere in any plane perpendicular to the z axis, we get an 8-cycle circular rainbow exactly like figure a. This is required because $L_z = 8\hbar$. But if we take a slice perpendicular to some other axis, such as the y axis, we don't get a circular rainbow as we would for a state with a definite value of L_y . It is obviously not possible to get circular rainbows for slices perpendicular to more than one axis.

For those with a taste for rigor, here is a complete argument:

Theorem: On the sphere, if a wavefunction has definite values of both L_z and L_x , then it is a wavefunction that is constant everywhere, so $\mathbf{L} = 0$.

Lemma 1: If the component of ℓ_A along a certain axis A has a definite value and is nonzero, then (a) $\Psi = 0$ at the poles, and (b) Ψ is of the form $Ae^{i\ell_A\phi}$ on any circle in a plane perpendicular to the axis. Part a holds because $\mathbf{L} = 0$ if $r_\perp = 0$. For b, see p. 921.

Lemma 2: If the component of \mathbf{L} along a certain axis has a definite value and is zero, then Ψ is constant in any plane perpendicular to that axis. This follows from lemma 1 in the case where $\ell_A = 0$.

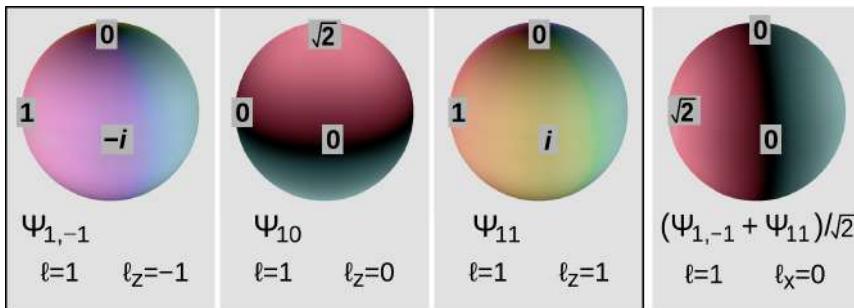
Case I: ℓ_z and ℓ_x are both nonzero. We have $\Psi = 0$ at the poles along both the x axis and the z axis. The z -axis pole is a point on the great circle perpendicular to the x axis, and vice versa, so applying 1b, $A = 0$ and Ψ vanishes on both of these great circles. But now if we apply 1b along any slice perpendicular to either axis, we get $\Psi = 0$ everywhere on that slice, so $\Psi = 0$ everywhere.

Case II: ℓ_z and ℓ_x are both zero. By lemma 2, Ψ is a constant everywhere.

Case III: One component is zero and the other nonzero. Let ℓ_z be the one that is zero. By 1a, $\Psi = 0$ at the x -axis pole, so by 2, $\Psi = 0$ on the great circle perpendicular to z . But then 1b tells us that $\Psi = 0$ everywhere.

13.4.3 Quantum numbers

Completeness



d / The three states inside the box are a complete set of quantum numbers for $\ell = 1$. Other states with $\ell = 1$, such as the one on the right, are not really new: they can be expressed as superpositions of the original three we chose.

For a given ℓ , consider the set of states with all the possible values of the angular momentum's component along some fixed axis. This set of states is *complete*, meaning that they encompass all the possible states with this ℓ .

For example, figure d shows wavefunctions with $\ell = 1$ that are solutions of the Schrödinger equation for a particle that is confined to the surface of a sphere. Although the formulae for these wavefunctions are not particularly complicated,⁸ they are not our main focus here, so to help with getting a feel for the idea of completeness, I have simply selected three points on the sphere at which to give numerical samples of the value of the wavefunction. These are the top (where the sphere is intersected by the positive z axis), left (x), and front (y). (Although the wavefunctions are shown using the color conventions defined in figure u, p. 916, these numerical samples should make the example understandable if you're looking at a black and white copy of the book.)

Suppose we arbitrarily choose the z axis as the one along which to quantize the component of the angular momentum. With this choice, we have three possible values for ℓ_z : -1 , 0 , and 1 . These three states are shown in the three boxes surrounded by the black rectangle. This set of three states is complete.

Consider, for example, the fourth state, shown on the right outside the box. This state is clearly identifiable as a copy of the $\ell_z = 0$ state, rotated by 90 degrees counterclockwise, so it is the $\ell_x = 0$ state. We might imagine that this would be an entirely new prize to be added to our stamp collection. But it is actually not a state that we didn't possess before. We can obtain it as the sum of the $\ell_z = -1$ and $\ell_z = 1$ states, divided by an appropriate normalization factor. Although I'm avoiding making this example an exercise in

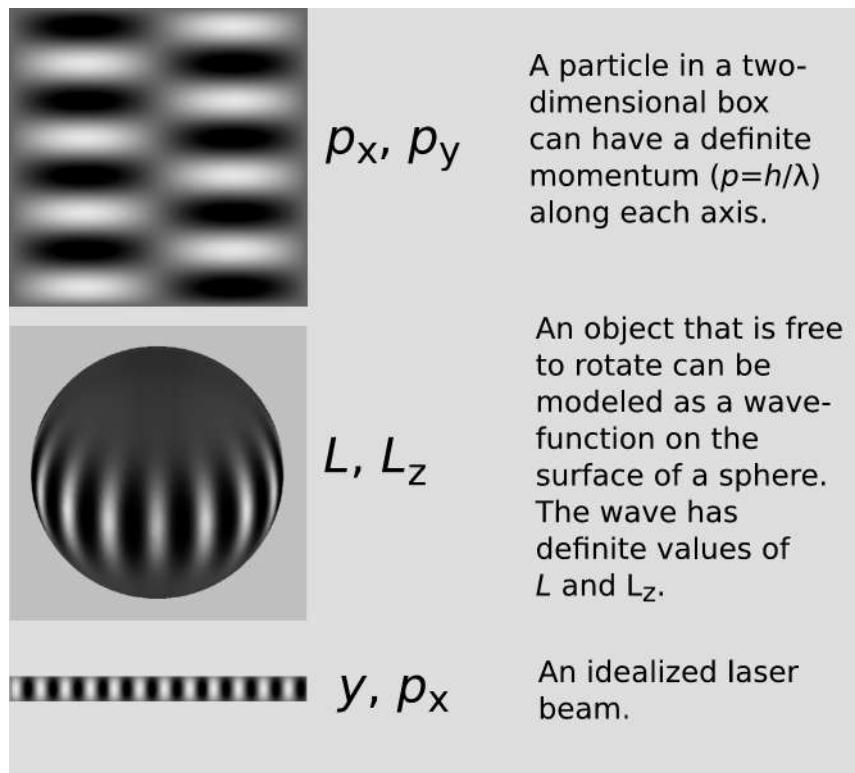
⁸They are $\Psi_{1,-1} = \sin \theta e^{-i\phi}$, $\Psi_{10} = \sqrt{2} \cos \theta$, and $\Psi_{11} = \sin \theta e^{i\phi}$, where θ is the angle measured down from the z axis, and ϕ is the angle running counterclockwise around the z axis. These functions are called spherical harmonics.

manipulating formulae, it is easy to check that the sum does work out properly at the three sample points.

Sets of compatible quantum numbers

Figure e shows some examples in which we can completely describe a wavefunction by giving a few numbers. These are referred to as “quantum numbers.” It is important that the quantum numbers we use in describing a state be compatible. By analogy, “Bond, James, 007” would be a clear and consistent definition of the famous fictional spy, but in general this identification scheme would not work, because although almost everyone has a first and last name, most people do not have a license to kill with a corresponding double-oh number.

e / Three examples of sets of compatible quantum numbers.



The laser beam in the figure is a state described according to its definite values p_x and y , so we have the vanishing uncertainties $\Delta p_x = 0$ and $\Delta y = 0$. Since the Heisenberg uncertainty principle doesn't talk about an x momentum in relation to a y position, this is OK. If we had been in doubt about whether this violated the uncertainty principle, we would have been reassured by our ability to draw the picture.

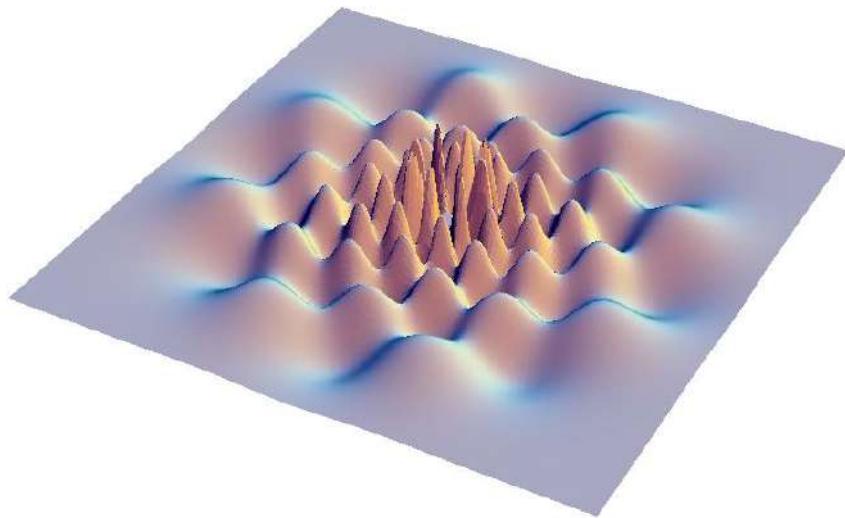
It is also possible to have *incompatible* quantum numbers. The combination of p_x with x would be an incompatible set of quantum numbers, because a state can't have a definite p_x and also a definite

x . If we try to draw such a wave, we fail. L_x and L_z would also be an incompatible set.

Complete and compatible sets of quantum numbers

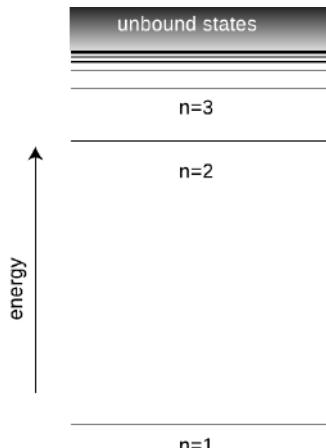
Let's summarize. Just as we expect everyone to have a first and last name, we expect there to be a complete and compatible set of quantum numbers for any given quantum-mechanical system. Completeness means that we have enough quantum numbers to uniquely describe every possible state of the system, although we may need to describe a state as a superposition, as with the state $\ell_x = 0$ in figure d on p. 925. Compatibility means that when we specify a set of quantum numbers, we aren't making a set of demands that can't be met. These ideas are revisited in a slightly fancier mathematical way on p. 990.

13.4.4 The hydrogen atom



f / A cross-section of a hydrogen wavefunction.

Deriving all the wavefunctions of the states of the hydrogen atom from first principles would be mathematically too complex for this book. (The ground state is not too hard, and we analyze it on p. 933.). But it's not hard to understand the logic behind the wavefunctions in visual terms. Consider the wavefunction from the beginning of the section, which is reproduced in figure f. Although the graph looks three-dimensional, it is really only a representation of the part of the wavefunction lying within a two-dimensional plane. The third (up-down) dimension of the plot represents the value of the wavefunction at a given point, not the third dimension of space. The plane chosen for the graph is the one perpendicular to the angular momentum vector.



g / The energy of a state in the hydrogen atom depends only on its n quantum number.

Each ring of peaks and valleys has eight wavelengths going around in a circle, so this state has $L = 8\hbar$, i.e., we label it $\ell = 8$. The wavelength is shorter near the center, and this makes sense because when the electron is close to the nucleus it has a lower electrical energy, a higher kinetic energy, and a higher momentum.

Between each ring of peaks in this wavefunction is a nodal circle, i.e., a circle on which the wavefunction is zero. The full three-dimensional wavefunction has nodal spheres: a series of nested spherical surfaces on which it is zero. The number of radii at which nodes occur, including $r = \infty$, is called n , and n turns out to be closely related to energy. The ground state has $n = 1$ (a single node only at $r = \infty$), and higher-energy states have higher n values. There is a simple equation relating n to energy, which we will discuss in subsection 13.4.5.

The numbers n and ℓ , which identify the state, are called its quantum numbers. A state of a given n and ℓ can be oriented in a variety of directions in space. We might try to indicate the orientation using the three quantum numbers $\ell_x = L_x/\hbar$, $\ell_y = L_y/\hbar$, and $\ell_z = L_z/\hbar$. But we have already seen that it is impossible to know all three of these simultaneously. To give the most complete possible description of a state, we choose an arbitrary axis, say the z axis, and label the state according to n , ℓ , and ℓ_z .⁹

Angular momentum requires motion, and motion implies kinetic energy. Thus it is not possible to have a given amount of angular momentum without having a certain amount of kinetic energy as well. Since energy relates to the n quantum number, this means that for a given n value there will be a maximum possible ℓ . It turns out that this maximum value of equals $n - 1$.

In general, we can list the possible combinations of quantum numbers as follows:

n can equal 1, 2, 3, ...
ℓ can range from 0 to $n - 1$, in steps of 1
ℓ_z can range from $-\ell$ to ℓ , in steps of 1

Applying these rules, we have the following list of states:

$n = 1, \ell = 0, \ell_z = 0$	one state
$n = 2, \ell = 0, \ell_z = 0$	one state
$n = 2, \ell = 1, \ell_z = -1, 0, \text{ or } 1$	three states
...	

self-check J

Continue the list for $n = 3$.

▷ Answer, p. 1067

Because the energy only depends on n , we have degeneracies. For example, the $n = 2$ energy level is 4-fold degenerate (and in fact

⁹See page 938 for a note about the two different systems of notations that are used for quantum numbers.

this degeneracy will be doubled to 8 when we take into account the intrinsic spin of the electron, sec. 13.4.6, p. 936). The degeneracy of the different ℓ_z states follows from symmetry, as in our original example of degeneracy on p. 922, and is therefore exact. The degeneracy with respect to different values of ℓ for the same n is not at all obvious, and is in fact not exact when effects such as relativity are taken into account. We refer to this as an “accidental” degeneracy. The very high level of degeneracy in the hydrogen atom means that when you observe it the hydrogen spectrum in your lab course, there is a great deal of structure that is effectively hidden from you. Historically, physicists were fooled by the apparent simplicity of the spectrum, and more than 70 years passed between the measurement of the spectrum and the time when the degeneracies were fully recognized and understood.

Figure h on page 930 shows the lowest-energy states of the hydrogen atom. The left-hand column of graphs displays the wavefunctions in the $x - y$ plane, and the right-hand column shows the probability distribution in a three-dimensional representation.

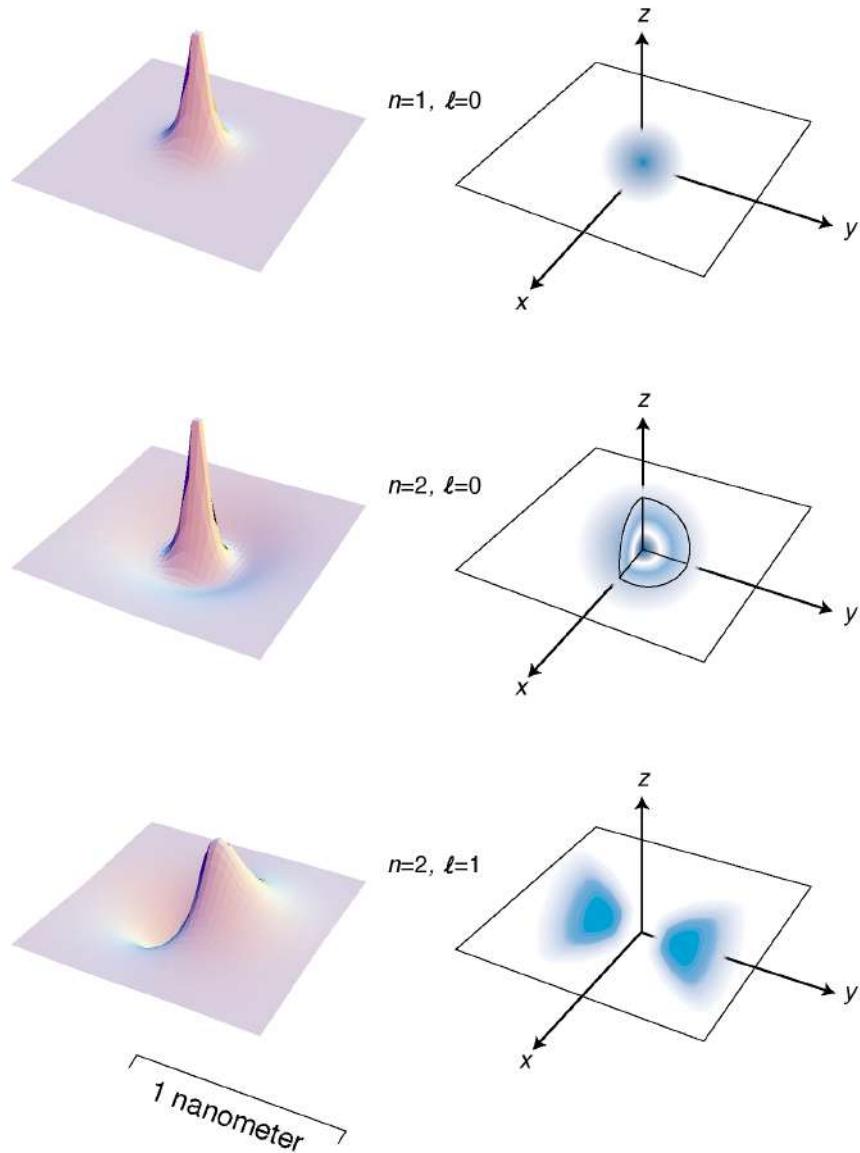
Discussion Questions

- A** The quantum number n is defined as the number of radii at which the wavefunction is zero, including $r = \infty$. Relate this to the features of figure h.
- B** Based on the definition of n , why can't there be any such thing as an $n = 0$ state?
- C** Relate the features of the wavefunction plots in figure h to the corresponding features of the probability distribution pictures.
- D** How can you tell from the wavefunction plots in figure h which ones have which angular momenta?
- E** Criticize the following incorrect statement: “The $\ell = 8$ wavefunction in figure f has a shorter wavelength in the center because in the center the electron is in a higher energy level.”
- F** Discuss the implications of the fact that the probability cloud in of the $n = 2, \ell = 1$ state is split into two parts.

13.4.5 Energies of states in hydrogen

History

The experimental technique for measuring the energy levels of an atom accurately is spectroscopy: the study of the spectrum of light emitted (or absorbed) by the atom. Only photons with certain energies can be emitted or absorbed by a hydrogen atom, for example, since the amount of energy gained or lost by the atom must equal the difference in energy between the atom's initial and final states. Spectroscopy had become a highly developed art several decades before Einstein even proposed the photon, and the Swiss



h / The three states of the hydrogen atom having the lowest energies.

spectroscopist Johann Balmer determined in 1885 that there was a simple equation that gave all the wavelengths emitted by hydrogen. In modern terms, we think of the photon wavelengths merely as indirect evidence about the underlying energy levels of the atom, and we rework Balmer's result into an equation for these atomic energy levels:

$$E_n = -\frac{2.2 \times 10^{-18} \text{ J}}{n^2},$$

This energy includes both the kinetic energy of the electron and the electrical energy. The zero-level of the electrical energy scale is chosen to be the energy of an electron and a proton that are

infinitely far apart. With this choice, negative energies correspond to bound states and positive energies to unbound ones.

Where does the mysterious numerical factor of 2.2×10^{-18} J come from? In 1913 the Danish theorist Niels Bohr realized that it was exactly numerically equal to a certain combination of fundamental physical constants:

$$E_n = -\frac{mk^2e^4}{2\hbar^2} \cdot \frac{1}{n^2},$$

where m is the mass of the electron, and k is the Coulomb force constant for electric forces.

Bohr was able to cook up a derivation of this equation based on the incomplete version of quantum physics that had been developed by that time, but his derivation is today mainly of historical interest. It assumes that the electron follows a circular path, whereas the whole concept of a path for a particle is considered meaningless in our more complete modern version of quantum physics. Although Bohr was able to produce the right equation for the energy levels, his model also gave various wrong results, such as predicting that the atom would be flat, and that the ground state would have $\ell = 1$ rather than the correct $\ell = 0$.

Approximate treatment

Rather than leaping straight into a full mathematical treatment, we'll start by looking for some physical insight, which will lead to an approximate argument that correctly reproduces the form of the Bohr equation.

A typical standing-wave pattern for the electron consists of a central oscillating area surrounded by a region in which the wavefunction tails off. As discussed in subsection 13.3.6, the oscillating type of pattern is typically encountered in the classically allowed region, while the tailing off occurs in the classically forbidden region where the electron has insufficient kinetic energy to penetrate according to classical physics. We use the symbol r for the radius of the spherical boundary between the classically allowed and classically forbidden regions. Classically, r would be the distance from the proton at which the electron would have to stop, turn around, and head back in.

If r had the same value for every standing-wave pattern, then we'd essentially be solving the particle-in-a-box problem in three dimensions, with the box being a spherical cavity. Consider the energy levels of the particle in a box compared to those of the hydrogen atom, i. They're qualitatively different. The energy levels of the particle in a box get farther and farther apart as we go higher in energy, and this feature doesn't even depend on the details of whether the box is two-dimensional or three-dimensional, or its exact shape. The reason for the spreading is that the box is taken to

particle in a box

hydrogen atom

i / The energy levels of a particle in a box, contrasted with those of the hydrogen atom.

be completely impenetrable, so its size, r , is fixed. A wave pattern with n humps has a wavelength proportional to r/n , and therefore a momentum proportional to n , and an energy proportional to n^2 . In the hydrogen atom, however, the force keeping the electron bound isn't an infinite force encountered when it bounces off of a wall, it's the attractive electrical force from the nucleus. If we put more energy into the electron, it's like throwing a ball upward with a higher energy — it will get farther out before coming back down. This means that in the hydrogen atom, we expect r to increase as we go to states of higher energy. This tends to keep the wavelengths of the high energy states from getting too short, reducing their kinetic energy. The closer and closer crowding of the energy levels in hydrogen also makes sense because we know that there is a certain energy that would be enough to make the electron escape completely, and therefore the sequence of bound states cannot extend above that energy.

When the electron is at the maximum classically allowed distance r from the proton, it has zero kinetic energy. Thus when the electron is at distance r , its energy is purely electrical:

$$[1] \quad E = -\frac{ke^2}{r}$$

Now comes the approximation. In reality, the electron's wavelength cannot be constant in the classically allowed region, but we pretend that it is. Since n is the number of nodes in the wavefunction, we can interpret it approximately as the number of wavelengths that fit across the diameter $2r$. We are not even attempting a derivation that would produce all the correct numerical factors like 2 and π and so on, so we simply make the approximation

$$[2] \quad \lambda \sim \frac{r}{n}.$$

Finally we assume that the typical kinetic energy of the electron is on the same order of magnitude as the absolute value of its total energy. (This is true to within a factor of two for a typical classical system like a planet in a circular orbit around the sun.) We then have

$$\begin{aligned} [3] \quad & \text{absolute value of total energy} \\ &= \frac{ke^2}{r} \\ &\sim K \\ &= p^2/2m \\ &= (h/\lambda)^2/2m \\ &\sim h^2 n^2 / 2mr^2 \end{aligned}$$

We now solve the equation $ke^2/r \sim h^2n^2/2mr^2$ for r and throw away numerical factors we can't hope to have gotten right, yielding

$$[4] \quad r \sim \frac{h^2 n^2}{m k e^2}.$$

Plugging $n = 1$ into this equation gives $r = 2 \text{ nm}$, which is indeed on the right order of magnitude. Finally we combine equations [4] and [1] to find

$$E \sim -\frac{m k^2 e^4}{h^2 n^2},$$

which is correct except for the numerical factors we never aimed to find.

Exact treatment of the ground state

The general proof of the Bohr equation for all values of n is beyond the mathematical scope of this book, but it's fairly straightforward to verify it for a particular n , especially given a lucky guess as to what functional form to try for the wavefunction. The form that works for the ground state is

$$\Psi = u e^{-r/a},$$

where $r = \sqrt{x^2 + y^2 + z^2}$ is the electron's distance from the proton, and u provides for normalization. In the following, the result $\partial r / \partial x = x/r$ comes in handy. Computing the partial derivatives that occur in the Laplacian, we obtain for the x term

$$\begin{aligned} \frac{\partial \Psi}{\partial x} &= \frac{\partial \Psi}{\partial r} \frac{\partial r}{\partial x} \\ &= -\frac{x}{ar} \Psi \\ \frac{\partial^2 \Psi}{\partial x^2} &= -\frac{1}{ar} \Psi - \frac{x}{a} \left(\frac{\partial}{\partial x} \frac{1}{r} \right) \Psi + \left(\frac{x}{ar} \right)^2 \Psi \\ &= -\frac{1}{ar} \Psi + \frac{x^2}{ar^3} \Psi + \left(\frac{x}{ar} \right)^2 \Psi, \end{aligned}$$

so

$$\nabla^2 \Psi = \left(-\frac{2}{ar} + \frac{1}{a^2} \right) \Psi.$$

The Schrödinger equation gives

$$\begin{aligned} E \cdot \Psi &= -\frac{\hbar^2}{2m} \nabla^2 \Psi + U \cdot \Psi \\ &= \frac{\hbar^2}{2m} \left(\frac{2}{ar} - \frac{1}{a^2} \right) \Psi - \frac{ke^2}{r} \cdot \Psi \end{aligned}$$

If we require this equation to hold for all r , then we must have equality for both the terms of the form $(\text{constant}) \times \Psi$ and for those

of the form $(\text{constant}/r) \times \Psi$. That means

$$E = -\frac{\hbar^2}{2ma^2}$$

and

$$0 = \frac{\hbar^2}{mar} - \frac{ke^2}{r}.$$

These two equations can be solved for the unknowns a and E , giving

$$a = \frac{\hbar^2}{mke^2}$$

and

$$E = -\frac{mk^2e^4}{2\hbar^2},$$

where the result for the energy agrees with the Bohr equation for $n = 1$. The calculation of the normalization constant u is relegated to homework problem 36.

self-check K

We've verified that the function $\Psi = he^{-r/a}$ is a solution to the Schrödinger equation, and yet it has a kink in it at $r = 0$. What's going on here? Didn't I argue before that kinks are unphysical? ▷ Answer, p. 1067

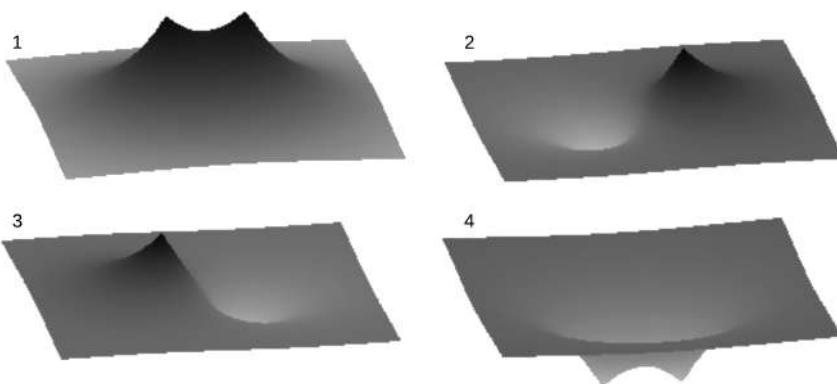
Wave phases in the hydrogen molecule

example 24

In example 15 on page 901, I argued that the existence of the H_2 molecule could essentially be explained by a particle-in-a-box argument: the molecule is a bigger box than an individual atom, so each electron's wavelength can be longer, its kinetic energy lower. Now that we're in possession of a mathematical expression for the wavefunction of the hydrogen atom in its ground state, we can make this argument a little more rigorous and detailed. Suppose that two hydrogen atoms are in a relatively cool sample of monoatomic hydrogen gas. Because the gas is cool, we can assume that the atoms are in their ground states. Now suppose that the two atoms approach one another. Making use again of the assumption that the gas is cool, it is reasonable to imagine that the atoms approach one another slowly. Now the atoms come a little closer, but still far enough apart that the region between them is classically forbidden. Each electron can tunnel through this classically forbidden region, but the tunneling probability is small. Each one is now found with, say, 99% probability in its original home, but with 1% probability in the other nucleus. Each electron is now in a state consisting of a superposition of the ground state of its own atom with the ground state of the other atom. There are two peaks in the superposed wavefunction, but one is a much bigger peak than the other.

An interesting question now arises. What are the relative phases of the two electrons? As discussed on page 895, the *absolute* phase of an electron's wavefunction is not really a meaningful concept. Suppose atom A contains electron Alice, and B electron Bob. Just before the collision, Alice may have wondered, "Is my phase positive right now, or is it negative? But of course I shouldn't ask myself such silly questions," she adds sheepishly.

j / Example 24.



But *relative* phases are well defined. As the two atoms draw closer and closer together, the tunneling probability rises, and eventually gets so high that each electron is spending essentially 50% of its time in each atom. It's now reasonable to imagine that either one of two possibilities could obtain. Alice's wavefunction could either look like j/1, with the two peaks in phase with one another, or it could look like j/2, with opposite phases. Because *relative* phases of wavefunctions are well defined, states 1 and 2 are physically distinguishable.¹⁰ In particular, the kinetic energy of state 2 is much higher; roughly speaking, it is like the two-hump wave pattern of the particle in a box, as opposed to 1, which looks roughly like the one-hump pattern with a much longer wavelength. Not only that, but an electron in state 1 has a large probability of being found in the central region, where it has a large negative electrical energy due to its interaction with both protons. State 2, on the other hand, has a low probability of existing in that region. Thus state 1 represents the true ground-state wavefunction of the H₂ molecule, and putting both Alice and Bob in that state results in a lower energy than their total energy when separated, so the molecule is bound, and will not fly apart spontaneously.

¹⁰The reader who has studied chemistry may find it helpful to make contact with the terminology and notation used by chemists. The state represented by pictures 1 and 4 is known as a σ orbital, which is a type of "bonding orbital." The state in 2 and 3 is a σ^* , a kind of "antibonding orbital." Note that although we will not discuss electron spin or the Pauli exclusion principle until sec. 13.4.6, p. 936, those considerations have no effect on this example, since the two electrons can have opposite spins.

State $j/3$, on the other hand, is not physically distinguishable from $j/2$, nor is $j/4$ from $j/1$. Alice may say to Bob, “Isn’t it wonderful that we’re in state 1 or 4? I love being stable like this.” But she knows it’s not meaningful to ask herself at a given moment which state she’s in, 1 or 4.

Discussion Questions

A States of hydrogen with n greater than about 10 are never observed in the sun. Why might this be?

B Sketch graphs of r and E versus n for the hydrogen atom, and compare with analogous graphs for the one-dimensional particle in a box.

13.4.6 Electron spin

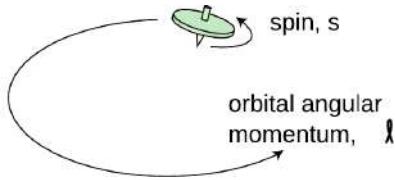
It’s disconcerting to the novice ping-pong player to encounter for the first time a more skilled player who can put spin on the ball. Even though you can’t see that the ball is spinning, you can tell something is going on by the way it interacts with other objects in its environment. In the same way, we can tell from the way electrons interact with other things that they have an intrinsic spin of their own. Experiments show that even when an electron is not moving through space, it still has angular momentum amounting to $\hbar/2$. An important historical experiment of this type, the Stern-Gerlach experiment, is described in detail in section 14.1.

This may seem paradoxical because the quantum moat, for instance, gave only angular momenta that were integer multiples of \hbar , not half-units, and I claimed that angular momentum was always quantized in units of \hbar , not just in the case of the quantum moat. That whole discussion, however, assumed that the angular momentum would come from the motion of a particle through space. The $\hbar/2$ angular momentum of the electron is simply a property of the particle, like its charge or its mass. It has nothing to do with whether the electron is moving or not, and it does not come from any internal motion within the electron. Nobody has ever succeeded in finding any internal structure inside the electron, and even if there was internal structure, it would be mathematically impossible for it to result in a half-unit of angular momentum.

We simply have to accept this $\hbar/2$ angular momentum, called the “spin” of the electron — Mother Nature rubs our noses in it as an observed fact. Protons and neutrons have the same $\hbar/2$ spin, while photons have an intrinsic spin of \hbar . In general, half-integer spins are typical of material particles. Integral values are found for the particles that carry forces: photons, which embody the electric and magnetic fields of force, as well as the more exotic messengers of the nuclear and gravitational forces. The photon is particularly important: it has spin 1.

As was the case with ordinary angular momentum, we can describe spin angular momentum in terms of its magnitude, and its

k / The top has angular momentum both because of the motion of its center of mass through space and due to its internal rotation. Electron spin is roughly analogous to the intrinsic spin of the top.



The diagram shows a top with a green circular base. A curved arrow around the top's axis represents orbital angular momentum. A small circle with a curved arrow inside represents spin. A vertical arrow pointing upwards represents the total angular momentum vector.

component along a given axis. We write s and s_z for these quantities, expressed in units of \hbar , so an electron has $s = 1/2$ and $s_z = +1/2$ or $-1/2$.

Odds and evens, and how they add up

From grade-school arithmetic, we have the rules

$$\text{even} + \text{even} = \text{even}$$

$$\text{odd} + \text{even} = \text{odd}$$

$$\text{odd} + \text{odd} = \text{even}.$$

Thus we know that $123456789 + 987654321$ is even, without having to actually compute the result. Dividing by two gives similar relationships for integer and half-integer angular momenta. For example, a half-integer plus an integer gives a half-integer, and therefore when we add the intrinsic spin $1/2$ of an electron to any additional, integer spin that the electron has from its motion through space, we get a half-integer angular momentum. That is, the *total* angular momentum of an electron will always be a half-integer. Similarly, when we add the intrinsic spin 1 of a photon to its angular momentum due to its integral motion through space, we will always get an integer. Thus the integer or half-integer character of any particle's *total* angular momentum (spin + motion) is determined entirely by the particle's spin.

These relationships tell us things about the spins we can make by putting together different particles to make bigger particles, and they also tell us things about decay processes.

Spin of the helium atom

example 25

A helium-4 atom consists of two protons, two neutrons, and two electrons. A proton, a neutron, and an electron each have spin $1/2$. Since the atom is a composite of six particles, each of which has half-integer spin, the atom as a whole has an integer angular momentum.

Emission of a photon from an atom

example 26

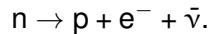
An atom can emit light,



This works in terms of angular momentum because the photon's spin 1 is an integer. Thus, regardless of whether the atom's angular momentum is an integer or a half-integer, the process is allowed by conservation of angular momentum. If the atom's angular momentum is an integer, then we have $\text{integer} = \text{integer} + 1$, and if it's a half-integer, $\text{half-integer} = \text{half-integer} + 1$; either of these is possible. If not for this logic, it would be impossible for matter to emit light. In general, if we want a particle such as a photon to pop into existence like this, it must have an integer spin.

Beta decay*example 27*

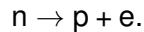
When a free neutron undergoes beta decay, we have



All four of these particles have spin 1/2, so the angular momenta go like

$$\text{half-integer} \rightarrow \text{half-integer} + \text{half-integer} + \text{half-integer},$$

which is possible, e.g., $1/2 = 3/2 - 5/2 + 3/2$. Because the neutrino has almost no interaction with normal matter, it normally flies off undetected, and the reaction was originally thought to be



With hindsight, this is impossible, because we can never have

$$\text{half-integer} \rightarrow \text{half-integer} + \text{half-integer}.$$

The reasoning holds not just for the beta decay of a free neutron, but for any beta decay: a neutrino or antineutrino must be emitted in order to conserve angular momentum. But historically, this was not understood at first, and when Enrico Fermi proposed the existence of the neutrino in 1934, the journal to which he first submitted his paper rejected it as “too remote from reality.”

States in hydrogen, with spin

Taking electron spin into account, we need a total of four quantum numbers to label a state of an electron in the hydrogen atom: n , ℓ , ℓ_z , and s_z . (We omit s because it always has the same value.) The symbols ℓ and ℓ_z include only the angular momentum the electron has because it is moving through space, not its spin angular momentum. The availability of two possible spin states of the electron leads to a doubling of the numbers of states:

$n = 1, \ell = 0, \ell_z = 0,$	$s_z = +1/2 \text{ or } -1/2$	two states
$n = 2, \ell = 0, \ell_z = 0,$	$s_z = +1/2 \text{ or } -1/2$	two states
$n = 2, \ell = 1, \ell_z = -1, 0, \text{ or } 1, s_z = +1/2 \text{ or } -1/2$		six states
...		

A note about notation

There are unfortunately two inconsistent systems of notation for the quantum numbers we've been discussing. The notation I've been using is the one that is used in nuclear physics, but there is a different one that is used in atomic physics.

nuclear physics	atomic physics
n	same
ℓ	same
ℓ_x	no notation
ℓ_y	no notation
ℓ_z	m
$s = 1/2$	no notation (sometimes σ)
s_x	no notation
s_y	no notation
s_z	s

The nuclear physics notation is more logical (not giving special status to the z axis) and more memorable (ℓ_z rather than the obscure m), which is why I use it consistently in this book, even though nearly all the applications we'll consider are atomic ones.

We are further encumbered with the following historically derived letter labels, which deserve to be eliminated in favor of the simpler numerical ones:

$\ell = 0$	$\ell = 1$	$\ell = 2$	$\ell = 3$			
s	p	d	f			
$n = 1$	$n = 2$	$n = 3$	$n = 4$	$n = 5$	$n = 6$	$n = 7$
K	L	M	N	O	P	Q

The spdf labels are used in both nuclear¹¹ and atomic physics, while the KLMNOPQ letters are used only to refer to states of electrons.

And finally, there is a piece of notation that is good and useful, but which I simply haven't mentioned yet. The vector $\mathbf{j} = \ell + \mathbf{s}$ stands for the total angular momentum of a particle in units of \hbar , including both orbital and spin parts. This quantum number turns out to be very useful in nuclear physics, because nuclear forces tend to exchange orbital and spin angular momentum, so a given energy level often contains a mixture of ℓ and s values, while remaining fairly pure in terms of j .

13.4.7 Atoms with more than one electron

What about other atoms besides hydrogen? It would seem that things would get much more complex with the addition of a second electron. A hydrogen atom only has one particle that moves around much, since the nucleus is so heavy and nearly immobile. Helium, with two, would be a mess. Instead of a wavefunction whose square tells us the probability of finding a single electron at any given location in space, a helium atom would need to have a wavefunction whose square would tell us the probability of finding two electrons at any given combination of points. Ouch! In addition, we would

¹¹After f, the series continues in alphabetical order. In nuclei that are spinning rapidly enough that they are almost breaking apart, individual protons and neutrons can be stirred up to ℓ values as high as 7, which is j.

have the extra complication of the electrical interaction between the two electrons, rather than being able to imagine everything in terms of an electron moving in a static field of force created by the nucleus alone.

Despite all this, it turns out that we can get a surprisingly good description of many-electron atoms simply by assuming the electrons can occupy the same standing-wave patterns that exist in a hydrogen atom. The ground state of helium, for example, would have both electrons in states that are very similar to the $n = 1$ states of hydrogen. The second-lowest-energy state of helium would have one electron in an $n = 1$ state, and the other in an $n = 2$ state. The relatively complex spectra of elements heavier than hydrogen can be understood as arising from the great number of possible combinations of states for the electrons.

A surprising thing happens, however, with lithium, the three-electron atom. We would expect the ground state of this atom to be one in which all three electrons settle down into $n = 1$ states. What really happens is that two electrons go into $n = 1$ states, but the third stays up in an $n = 2$ state. This is a consequence of a new principle of physics:

The Pauli Exclusion Principle

Two electrons can never occupy the same state.¹²

There are two $n = 1$ states, one with $s_z = +1/2$ and one with $s_z = -1/2$, but there is no third $n = 1$ state for lithium's third electron to occupy, so it is forced to go into an $n = 2$ state.

It can be proved mathematically that the Pauli exclusion principle applies to any type of particle that has half-integer spin. Thus two neutrons can never occupy the same state, and likewise for two protons. Photons, however, are immune to the exclusion principle because their spin is an integer.

Deriving the periodic table

We can now account for the structure of the periodic table, which seemed so mysterious even to its inventor Mendeleev. The first row consists of atoms with electrons only in the $n = 1$ states:

- H 1 electron in an $n = 1$ state
- He 2 electrons in the two $n = 1$ states

The next row is built by filling the $n = 2$ energy levels:

- Li 2 electrons in $n = 1$ states, 1 electron in an $n = 2$ state
- Be 2 electrons in $n = 1$ states, 2 electrons in $n = 2$ states
- ...
- O 2 electrons in $n = 1$ states, 6 electrons in $n = 2$ states
- F 2 electrons in $n = 1$ states, 7 electrons in $n = 2$ states
- Ne 2 electrons in $n = 1$ states, 8 electrons in $n = 2$ states

In the third row we start in on the $n = 3$ levels:

Na 2 electrons in $n = 1$ states, 8 electrons in $n = 2$ states, 1 electron in an $n = 3$ state

...

We can now see a logical link between the filling of the energy levels and the structure of the periodic table. Column 0, for example, consists of atoms with the right number of electrons to fill all the available states up to a certain value of n . Column I contains atoms like lithium that have just one electron more than that.

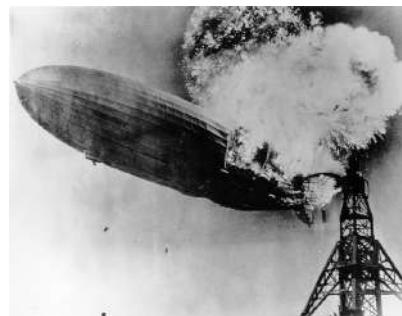
This shows that the columns relate to the filling of energy levels, but why does that have anything to do with chemistry? Why, for example, are the elements in columns I and VII dangerously reactive? Consider, for example, the element sodium (Na), which is so reactive that it may burst into flames when exposed to air. The electron in the $n = 3$ state has an unusually high energy. If we let a sodium atom come in contact with an oxygen atom, energy can be released by transferring the $n = 3$ electron from the sodium to one of the vacant lower-energy $n = 2$ states in the oxygen. This energy is transformed into heat. Any atom in column I is highly reactive for the same reason: it can release energy by giving away the electron that has an unusually high energy.

Column VII is spectacularly reactive for the opposite reason: these atoms have a single vacancy in a low-energy state, so energy is released when these atoms steal an electron from another atom.

It might seem as though these arguments would only explain reactions of atoms that are in different rows of the periodic table, because only in these reactions can a transferred electron move from a higher- n state to a lower- n state. This is incorrect. An $n = 2$ electron in fluorine (F), for example, would have a different energy than an $n = 2$ electron in lithium (Li), due to the different number of protons and electrons with which it is interacting. Roughly speaking, the $n = 2$ electron in fluorine is more tightly bound (lower in energy) because of the larger number of protons attracting it. The effect of the increased number of attracting protons is only partly counteracted by the increase in the number of repelling electrons, because the forces exerted on an electron by the other electrons are in many different directions and cancel out partially.

I	II	III	IV	V	VI	VII	0
1 H							2 He
3 Li	4 Be	5 B	6 C	7 N	8 O	9 F	10 Ne
11 Na	...						

I / The beginning of the periodic table.



m / Hydrogen is highly reactive.

Problems

The symbols \checkmark , \blacksquare , etc. are explained on page 953.

- 1** If a radioactive substance has a half-life of one year, does this mean that it will be completely decayed after two years? Explain. \blacksquare

- 2** What is the probability of rolling a pair of dice and getting “snake eyes,” i.e., both dice come up with ones? \blacksquare

- 3** This problem has been deleted. \blacksquare

- 4** This problem has been deleted. \blacksquare

- 5** Refer to the probability distribution for people’s heights in figure f on page 864.

(a) Show that the graph is properly normalized.

(b) Estimate the fraction of the population having heights between 140 and 150 cm. \checkmark \blacksquare

- 6** (a) A nuclear physicist is studying a nuclear reaction caused in an accelerator experiment, with a beam of ions from the accelerator striking a thin metal foil and causing nuclear reactions when a nucleus from one of the beam ions happens to hit one of the nuclei in the target. After the experiment has been running for a few hours, a few billion radioactive atoms have been produced, embedded in the target. She does not know what nuclei are being produced, but she suspects they are an isotope of some heavy element such as Pb, Bi, Fr or U. Following one such experiment, she takes the target foil out of the accelerator, sticks it in front of a detector, measures the activity every 5 min, and makes a graph (figure). The isotopes she thinks may have been produced are:

isotope	half-life (minutes)
^{211}Pb	36.1
^{214}Pb	26.8
^{214}Bi	19.7
^{223}Fr	21.8
^{239}U	23.5

Which one is it?

- (b) Having decided that the original experimental conditions produced one specific isotope, she now tries using beams of ions traveling at several different speeds, which may cause different reactions. The following table gives the activity of the target 10, 20 and 30 minutes after the end of the experiment, for three different ion speeds.

	activity (millions of decays/s) after...		
	10 min	20 min	30 min
first ion speed	1.933	0.832	0.382
second ion speed	1.200	0.545	0.248
third ion speed	7.211	1.296	0.248

Since such a large number of decays is being counted, assume that

Problem 6.

the data are only inaccurate due to rounding off when writing down the table. Which are consistent with the production of a single isotope, and which imply that more than one isotope was being created? ■

7 Devise a method for testing experimentally the hypothesis that a gambler's chance of winning at craps is independent of her previous record of wins and losses. If you don't invoke the mathematical definition of statistical independence, then you haven't proposed a test. This has nothing to do with the details of the rules of craps, or with the fact that it's a game played using dice. ■

8 A blindfolded person fires a gun at a circular target of radius b , and is allowed to continue firing until a shot actually hits it. Any part of the target is equally likely to get hit. We measure the random distance r from the center of the circle to where the bullet went in.
(a) Show that the probability distribution of r must be of the form $D(r) = kr$, where k is some constant. (Of course we have $D(r) = 0$ for $r > b$.)

- (b) Determine k by requiring D to be properly normalized. ✓
(c) Find the average value of r . ✓
(d) Interpreting your result from part c, how does it compare with $b/2$? Does this make sense? Explain. ■

9 We are given some atoms of a certain radioactive isotope, with half-life $t_{1/2}$. We pick one atom at random, and observe it for one half-life, starting at time zero. If it decays during that one-half-life period, we record the time t at which the decay occurred. If it doesn't, we reset our clock to zero and keep trying until we get an atom that cooperates. The final result is a time $0 \leq t \leq t_{1/2}$, with a distribution that looks like the usual exponential decay curve, but with its tail chopped off.

- (a) Find the distribution $D(t)$, with the proper normalization. ✓
(b) Find the average value of t . ✓
(c) Interpreting your result from part b, how does it compare with $t_{1/2}/2$? Does this make sense? Explain. ■

10 The speed, v , of an atom in an ideal gas has a probability distribution of the form $D(v) = bve^{-cv^2}$, where $0 \leq v < \infty$, c relates to the temperature, and b is determined by normalization.

- (a) Sketch the distribution.
(b) Find b in terms of c . ✓
(c) Find the average speed in terms of c , eliminating b . (Don't try to do the indefinite integral, because it can't be done in closed form. The relevant definite integral can be found in tables or done with computer software.) ✓ ■

11 All helium on earth is from the decay of naturally occurring heavy radioactive elements such as uranium. Each alpha particle that is emitted ends up claiming two electrons, which makes it a helium atom. If the original ^{238}U atom is in solid rock (as opposed to the earth's molten regions), the He atoms are unable to diffuse out of the rock. This problem involves dating a rock using the known decay properties of uranium 238. Suppose a geologist finds a sample of hardened lava, melts it in a furnace, and finds that it contains 1230 mg of uranium and 2.3 mg of helium. ^{238}U decays by alpha emission, with a half-life of 4.5×10^9 years. The subsequent chain of alpha and electron (beta) decays involves much shorter half-lives, and terminates in the stable nucleus ^{206}Pb . Almost all natural uranium is ^{238}U , and the chemical composition of this rock indicates that there were no decay chains involved other than that of ^{238}U .

(a) How many alphas are emitted per decay chain? [Hint: Use conservation of mass.]

(b) How many electrons are emitted per decay chain? [Hint: Use conservation of charge.]

(c) How long has it been since the lava originally hardened? ✓ ■

12 When light is reflected from a mirror, perhaps only 80% of the energy comes back. One could try to explain this in two different ways: (1) 80% of the photons are reflected, or (2) all the photons are reflected, but each loses 20% of its energy. Based on your everyday knowledge about mirrors, how can you tell which interpretation is correct? [Based on a problem from PSSC Physics.] ■

13 Suppose we want to build an electronic light sensor using an apparatus like the one described in subsection 13.2.2 on p. 875. How would its ability to detect different parts of the spectrum depend on the type of metal used in the capacitor plates? ■

14 The photoelectric effect can occur not just for metal cathodes but for any substance, including living tissue. Ionization of DNA molecules can cause cancer or birth defects. If the energy required to ionize DNA is on the same order of magnitude as the energy required to produce the photoelectric effect in a metal, which of the following types of electromagnetic waves might pose such a hazard? Explain.

60 Hz waves from power lines

100 MHz FM radio

microwaves from a microwave oven

visible light

ultraviolet light

x-rays



15 (a) Rank-order the photons according to their wavelengths, frequencies, and energies. If two are equal, say so. Explain all your answers.

(b) Photon 3 was emitted by a xenon atom going from its second-lowest-energy state to its lowest-energy state. Which of photons 1, 2, and 4 are capable of exciting a xenon atom from its lowest-energy state to its second-lowest-energy state? Explain. ■

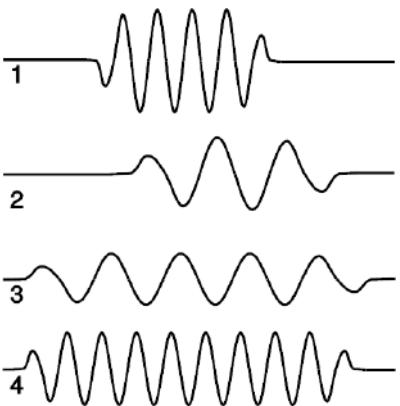
16 The figures show the wavefunction of an electron as a function of position. Which one could represent an electron speeding up as it moves to the right? Explain. ■

17 The beam of a 100 W overhead projector covers an area of $1 \text{ m} \times 1 \text{ m}$ when it hits the screen 3 m away. Estimate the number of photons that are in flight at any given time. (Since this is only an estimate, we can ignore the fact that the beam is not parallel.) ✓ ■

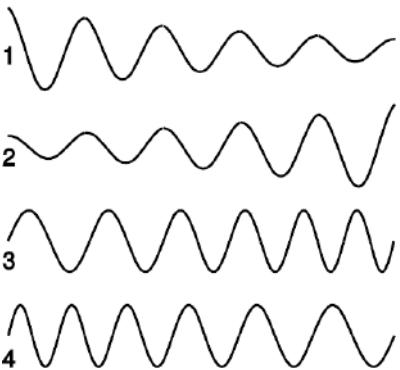
18 In the photoelectric effect, electrons are observed with virtually no time delay ($\sim 10 \text{ ns}$), even when the light source is very weak. (A weak light source does however only produce a small number of ejected electrons.) The purpose of this problem is to show that the lack of a significant time delay contradicted the classical wave theory of light, so throughout this problem you should put yourself in the shoes of a classical physicist and pretend you don't know about photons at all. At that time, it was thought that the electron might have a radius on the order of 10^{-15} m . (Recent experiments have shown that if the electron has any finite size at all, it is far smaller.)

(a) Estimate the power that would be soaked up by a single electron in a beam of light with an intensity of 1 mW/m^2 . ✓

(b) The energy, E_s , required for the electron to escape through the surface of the cathode is on the order of 10^{-19} J . Find how long it would take the electron to absorb this amount of energy, and explain why your result constitutes strong evidence that there is something wrong with the classical theory. ✓ ■



Problem 15.



Problem 16.

19 In a television, suppose the electrons are accelerated from rest through a voltage difference of 10^4 V. What is their final wavelength?

✓ ■

20 Use the Heisenberg uncertainty principle to estimate the minimum velocity of a proton or neutron in a ^{208}Pb nucleus, which has a diameter of about 13 fm ($1 \text{ fm} = 10^{-15} \text{ m}$). Assume that the speed is nonrelativistic, and then check at the end whether this assumption was warranted.

✓ ■

21 Find the energy of a nonrelativistic particle in a one-dimensional box of length L , expressing your result in terms of L , the particle's mass m , the number of peaks and valleys n in the wavefunction, and fundamental constants.

✓ ■

22 A free electron that contributes to the current in an ohmic material typically has a speed of 10^5 m/s (much greater than the drift velocity).

(a) Estimate its de Broglie wavelength, in nm. ✓

(b) If a computer memory chip contains 10^8 electric circuits in a 1 cm^2 area, estimate the linear size, in nm, of one such circuit. ✓

(c) Based on your answers from parts a and b, does an electrical engineer designing such a chip need to worry about wave effects such as diffraction?

(d) Estimate the maximum number of electric circuits that can fit on a 1 cm^2 computer chip before quantum-mechanical effects become important.

■

23 In classical mechanics, an interaction energy of the form $U(x) = \frac{1}{2}kx^2$ gives a harmonic oscillator: the particle moves back and forth at a frequency $\omega = \sqrt{k/m}$. This form for $U(x)$ is often a good approximation for an individual atom in a solid, which can vibrate around its equilibrium position at $x = 0$. (For simplicity, we restrict our treatment to one dimension, and we treat the atom as a single particle rather than as a nucleus surrounded by electrons). The atom, however, should be treated quantum-mechanically, not classically. It will have a wave function. We expect this wave function to have one or more peaks in the classically allowed region, and we expect it to tail off in the classically forbidden regions to the right and left. Since the shape of $U(x)$ is a parabola, not a series of flat steps as in figure m on page 908, the wavy part in the middle will not be a sine wave, and the tails will not be exponentials.

(a) Show that there is a solution to the Schrödinger equation of the form

$$\Psi(x) = e^{-bx^2},$$

and relate b to k , m , and \hbar . To do this, calculate the second derivative, plug the result into the Schrödinger equation, and then find what value of b would make the equation valid for *all* values of x . This wavefunction turns out to be the ground state. Note that this wavefunction is not properly normalized — don't worry about that.

- (b) Sketch a graph showing what this wavefunction looks like.
 (c) Let's interpret b . If you changed b , how would the wavefunction look different? Demonstrate by sketching two graphs, one for a smaller value of b , and one for a larger value.
 (d) Making k greater means making the atom more tightly bound. Mathematically, what happens to the value of b in your result from part a if you make k greater? Does this make sense physically when you compare with part c?

✓ ■

24 (a) A distance scale is shown below the wavefunctions and probability densities illustrated in figure h on page 930. Compare this with the order-of-magnitude estimate derived in section 13.4.5, p. 929, for the radius r at which the wavefunction begins tailing off. Was the estimate on the right order of magnitude?

(b) Although we normally say the moon orbits the earth, actually they both orbit around their common center of mass, which is below the earth's surface but not at its center. The same is true of the hydrogen atom. Does the center of mass lie inside the proton, or outside it?

■

25 The figure shows eight of the possible ways in which an electron in a hydrogen atom could drop from a higher energy state to a state of lower energy, releasing the difference in energy as a photon. Of these eight transitions, only D, E, and F produce photons with wavelengths in the visible spectrum.

(a) Which of the visible transitions would be closest to the violet end of the spectrum, and which would be closest to the red end? Explain.

(b) In what part of the electromagnetic spectrum would the photons from transitions A, B, and C lie? What about G and H? Explain.

(c) Is there an upper limit to the wavelengths that could be emitted by a hydrogen atom going from one bound state to another bound state? Is there a lower limit? Explain.

■

26 Find an equation for the wavelength of the photon emitted when the electron in a hydrogen atom makes a transition from energy level n_1 to level n_2 .

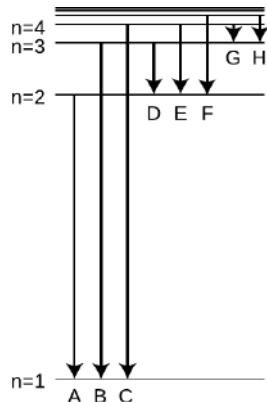
✓ ■

27 Estimate the angular momentum of a spinning basketball, in units of \hbar . Explain how this result relates to the correspondence principle.

■

28 Assume that the kinetic energy of an electron in the $n = 1$ state of a hydrogen atom is on the same order of magnitude as the absolute value of its total energy, and estimate a typical speed at which it would be moving. (It cannot really have a single, definite speed, because its kinetic and interaction energy trade off at different distances from the proton, but this is just a rough estimate of a typical speed.) Based on this speed, were we justified in assuming that the electron could be described nonrelativistically?

■



Problem 25.

29 Before the quantum theory, experimentalists noted that in many cases, they would find three lines in the spectrum of the same atom that satisfied the following mysterious rule: $1/\lambda_1 = 1/\lambda_2 + 1/\lambda_3$. Explain why this would occur. Do not use reasoning that only works for hydrogen — such combinations occur in the spectra of all elements. [Hint: Restate the equation in terms of the energies of photons.]

30 The wavefunction of the electron in the ground state of a hydrogen atom, shown in the top left of figure h on p. 930, is

$$\Psi = \pi^{-1/2} a^{-3/2} e^{-r/a},$$

where r is the distance from the proton, and $a = \hbar^2/kme^2 = 5.3 \times 10^{-11}$ m is a constant that sets the size of the wave. The figure doesn't show the proton; let's take the proton to be a sphere with a radius of $b = 0.5$ fm.

- (a) Reproduce figure h in a rough sketch, and indicate, relative to the size of your sketch, some idea of how big a and b are.
- (b) Calculate symbolically, without plugging in numbers, the probability that at any moment, the electron is inside the proton. [Hint: Does it matter if you plug in $r = 0$ or $r = b$ in the equation for the wavefunction?]
- (c) Calculate the probability numerically.
- (d) Based on the equation for the wavefunction, is it valid to think of a hydrogen atom as having a finite size? Can a be interpreted as the size of the atom, beyond which there is nothing? Or is there any limit on how far the electron can be from the proton?

31 Use physical reasoning to explain how the equation for the energy levels of hydrogen,

$$E_n = -\frac{mk^2e^4}{2\hbar^2} \cdot \frac{1}{n^2},$$

should be generalized to the case of an atom with atomic number Z that has had all its electrons removed except for one.

32 A muon is a subatomic particle that acts exactly like an electron except that its mass is 207 times greater. Muons can be created by cosmic rays, and it can happen that one of an atom's electrons is displaced by a muon, forming a muonic atom. If this happens to a hydrogen atom, the resulting system consists simply of a proton plus a muon.

- (a) Based on the results of section 13.4.5, how would the size of a muonic hydrogen atom in its ground state compare with the size of the normal atom?
- (b) If you were searching for muonic atoms in the sun or in the earth's atmosphere by spectroscopy, in what part of the electromagnetic spectrum would you expect to find the absorption lines?

33 An electron is initially at rest. A photon collides with the electron and rebounds from the collision at 180 degrees, i.e., going back along the path on which it came. The rebounding photon has a different energy, and therefore a different frequency and wavelength. Show that, based on conservation of energy and momentum, the difference between the photon's initial and final wavelengths must be $2h/mc$, where m is the mass of the electron. The experimental verification of this type of "pool-ball" behavior by Arthur Compton in 1923 was taken as definitive proof of the particle nature of light. Note that we're not making any nonrelativistic approximations. To keep the algebra simple, you should use natural units — in fact, it's a good idea to use even-more-natural-than-natural units, in which we have not just $c = 1$ but also $h = 1$, and $m = 1$ for the mass of the electron. You'll also probably want to use the relativistic relationship $E^2 - p^2 = m^2$, which becomes $E^2 - p^2 = 1$ for the energy and momentum of the electron in these units. ■

34 Generalize the result of problem 33 to the case where the photon bounces off at an angle other than 180° with respect to its initial direction of motion. ■

35 On page 908 we derived an expression for the probability that a particle would tunnel through a rectangular barrier, i.e., a region in which the interaction energy $U(x)$ has a graph that looks like a rectangle. Generalize this to a barrier of any shape. [Hints: First try generalizing to two rectangular barriers in a row, and then use a series of rectangular barriers to approximate the actual curve of an arbitrary function $U(x)$. Note that the width and height of the barrier in the original equation occur in such a way that all that matters is the area under the U -versus- x curve. Show that this is still true for a series of rectangular barriers, and generalize using an integral.] If you had done this calculation in the 1930's you could have become a famous physicist. ■

36 Show that the wavefunction given in problem 30 is properly normalized. ■

37 Show that a wavefunction of the form $\Psi = e^{by} \sin ax$ is a possible solution of the Schrödinger equation in two dimensions, with a constant potential U . Can we tell whether it would apply to a classically allowed region, or a classically forbidden one? ■

38 This problem generalizes the one-dimensional result from problem 21.

Find the energy levels of a particle in a three-dimensional rectangular box with sides of length a , b , and c . \checkmark ■

39 Americium-241 is an artificial isotope used in smoke detectors. It undergoes alpha decay, with a half-life of 432 years. As discussed in example 17 on page 909, alpha decay can be understood as a tunneling process, and although the barrier is not rectangular in shape, the equation for the tunneling probability on page 909 can still be used as a rough guide to our thinking. For americium-241, the tunneling probability is about 1×10^{-29} . Suppose that this nucleus were to decay by emitting a helium-3 nucleus instead of an alpha particle (helium-4). Estimate the relevant tunneling probability, assuming that the total energy E remains the same. This higher probability is contrary to the empirical observation that this nucleus is not observed to decay by ${}^3\text{He}$ emission with any significant probability, and in general ${}^3\text{He}$ emission is almost unknown in nature; this is mainly because the ${}^3\text{He}$ nucleus is far less stable than the helium-4 nucleus, and the difference in binding energy reduces the energy available for the decay. □

40 As far as we know, the mass of the photon is zero. However, it's not possible to prove by experiments that anything is zero; all we can do is put an upper limit on the number. As of 2008, the best experimental upper limit on the mass of the photon is about 1×10^{-52} kg. Suppose that the photon's mass really isn't zero, and that the value is at the top of the range that is consistent with the present experimental evidence. In this case, the c occurring in relativity would no longer be interpreted as the speed of light. As with material particles, the speed v of a photon would depend on its energy, and could never be as great as c . Estimate the relative size $(c - v)/c$ of the discrepancy in speed, in the case of a photon of visible light.

▷ Answer, p. 1069 □

41 Hydrogen is the only element whose energy levels can be expressed exactly in an equation. Calculate the ratio λ_E/λ_F of the wavelengths of the transitions labeled E and F in problem 25 on p. 947. Express your answer as an exact fraction, not a decimal approximation. In an experiment in which atomic wavelengths are being measured, this ratio provides a natural, stringent check on the precision of the results. ✓ □

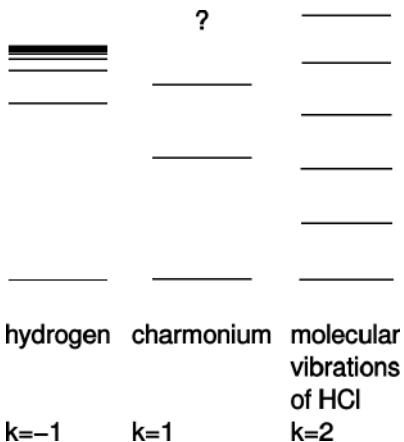
42 Give a numerical comparison of the number of photons per second emitted by a hundred-watt FM radio transmitter and a hundred-watt lightbulb. ✓ □

43 On pp. 931-933 of subsection 13.4.5, we used simple algebra to derive an approximate expression for the energies of states in hydrogen, without having to explicitly solve the Schrödinger equation. As input to the calculation, we used the the proportionality $U \propto r^{-1}$, which is a characteristic of the electrical interaction. The result for the energy of the n th standing wave pattern was $E_n \propto n^{-2}$.

There are other systems of physical interest in which we have $U \propto r^k$ for values of k besides -1 . Problem 23 discusses the ground state of the harmonic oscillator, with $k = 2$ (and a positive constant of proportionality). In particle physics, systems called charmonium and bottomonium are made out of pairs of subatomic particles called quarks, which interact according to $k = 1$, i.e., a force that is independent of distance. (Here we have a positive constant of proportionality, and $r > 0$ by definition. The motion turns out not to be too relativistic, so the Schrödinger equation is a reasonable approximation.) The figure shows actual energy levels for these three systems, drawn with different energy scales so that they can all be shown side by side. The sequence of energies in hydrogen approaches a limit, which is the energy required to ionize the atom. In charmonium, only the first three levels are known.¹³

Generalize the method used for $k = -1$ to any value of k , and find the exponent j in the resulting proportionality $E_n \propto n^j$. Compare the theoretical calculation with the behavior of the actual energies shown in the figure. Comment on the limit $k \rightarrow \infty$. ✓ ■

44 The electron, proton, and neutron were discovered, respectively, in 1897, 1919, and 1932. The neutron was late to the party, and some physicists felt that it was unnecessary to consider it as fundamental. Maybe it could be explained as simply a proton with an electron trapped inside it. The charges would cancel out, giving the composite particle the correct neutral charge, and the masses at least approximately made sense (a neutron is heavier than a proton). (a) Given that the diameter of a proton is on the order of 10^{-15} m, use the Heisenberg uncertainty principle to estimate the trapped electron's minimum momentum. ✓
 (b) Find the electron's minimum kinetic energy. ✓
 (c) Show via $E = mc^2$ that the proposed explanation may have a problem, because the contribution to the neutron's mass from the electron's kinetic energy would be comparable to the neutron's entire mass. ■



Problem 43.

¹³See Barnes et al., “The XYZs of Charmonium at BES,” arxiv.org/abs/hep-ph/0608103. To avoid complication, the levels shown are only those in the group known for historical reasons as the Ψ and J/Ψ .

45 Suppose that an electron, in one dimension, is confined to a certain region of space so that its wavefunction is given by

$$\Psi = \begin{cases} 0 & \text{if } x < 0 \\ A \sin(2\pi x/L) & \text{if } 0 \leq x \leq L \\ 0 & \text{if } x > L \end{cases}$$

Determine the constant A from normalization. ✓ ■

46 In the following, x and y are variables, while u and v are constants. Compute (a) $\partial(ux \ln(vy))/\partial x$, (b) $\partial(ux \ln(vy))/\partial y$. ✓ ■

47 (a) A radio transmitter radiates power P in all directions, so that the energy spreads out spherically. Find the energy density at a distance r . ✓

(b) Let the wavelength be λ . As described in example 8 on p. 878, find the number of photons in a volume λ^3 at this distance r . ✓

(c) For a 1000 kHz AM radio transmitting station, assuming reasonable values of P and r , verify, as claimed in the example, that the result from part b is very large. ■

48 The wavefunction Ψ of an electron is a complex number. Make up an example of a value for the wavefunction that is not a real number, and consider the following expressions: Ψ^2 , $|\Psi|^2$, $|\Psi^2|$. Which of these would it make sense to interpret as a probability density? All of them? Some? Only one? ▷ Solution, p. 1056 ■

49 (a) Consider the function defined by $f(x, y) = (x - y)^2$. Visualize the graph of this function as a surface. (This is a simple enough example that you should not have to resort to computer software.) Use this visualization to determine the behavior of the sign of the Laplacian, as in example 20 on p. 912.

(b) Consider the following incorrect calculation of this Laplacian. We take the first derivatives and find

$$\frac{\partial f}{\partial x} + \frac{\partial f}{\partial y} = 0.$$

Next we take the second derivatives, but those are zero as well, so the Laplacian is zero. Critique this calculation in two ways: (1) by comparing with part a; (2) by comparing with a correct calculation. (c) In general, if we have a function f of two variables, the quantity $Q = \partial f/\partial x + \partial f/\partial y$ can never be of physical interest, because it is not rotationally invariant (sec. 3.4.2, p. 195). Prove this by showing that by rotating your coordinate system, you can get a completely different answer than the one calculated in part b.

▷ Solution, p. 1056 ■

- 50** Let $\Psi = e^{2x+y}$. Compute $\nabla^2\Psi$. [If you get $9e^{2x+y}$, then you've made the mistake described in problem 49.] \checkmark ■

Key to symbols:

■ easy ■ typical ■ challenging ■ difficult ■ very difficult

\checkmark An answer check is available at www.lightandmatter.com.

Exercises

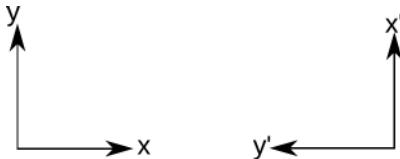
Exercise 13A: Quantum Versus Classical Randomness

1. Imagine the classical version of the particle in a one-dimensional box. Suppose you insert the particle in the box and give it a known, predetermined energy, but a random initial position and a random direction of motion. You then pick a random later moment in time to see where it is. Sketch the resulting probability distribution by shading on top of a line segment. Does the probability distribution depend on energy?
2. Do similar sketches for the first few energy levels of the quantum mechanical particle in a box, and compare with 1.
3. Do the same thing as in 1, but for a classical hydrogen atom in two dimensions, which acts just like a miniature solar system. Assume you're always starting out with the same fixed values of energy and angular momentum, but a position and direction of motion that are otherwise random. Do this for $L = 0$, and compare with a real $L = 0$ probability distribution for the hydrogen atom.
4. Repeat 3 for a nonzero value of L , say $L=\hbar$.
5. Summarize: Are the classical probability distributions accurate? What qualitative features are possessed by the classical diagrams but not by the quantum mechanical ones, or vice-versa?

Exercise 13B: Choice of quantum numbers

We could choose to identify a human by their first and last names, or by their social security number. We're free to choose any set of labels, as long as they're compatible.

1. Warm-up, nothing to do with quantum mechanics:



Express the x', y' coordinates in terms of the x, y coordinates:

$$x' = \boxed{} x + \boxed{} y$$
$$y' = \boxed{} x + \boxed{} y$$

2. Consider these four $\ell = 1$ wavefunctions on the “quantum moat:”

$$\Psi_{\circlearrowleft} = \boxed{} e^{i\theta} \quad \Psi_{\circlearrowright} = \boxed{} e^{-i\theta}$$
$$\Psi_s = \boxed{} \sin \theta \quad \Psi_c = \boxed{} \cos \theta$$

Determine the normalization factor for your group's wave.

3. We now want to discuss the standing waves in terms of the traveling waves:

$$\Psi_c = \boxed{} \Psi_{\circlearrowleft} + \boxed{} \Psi_{\circlearrowright}$$
$$\Psi_s = \boxed{} \Psi_{\circlearrowleft} + \boxed{} \Psi_{\circlearrowright}.$$

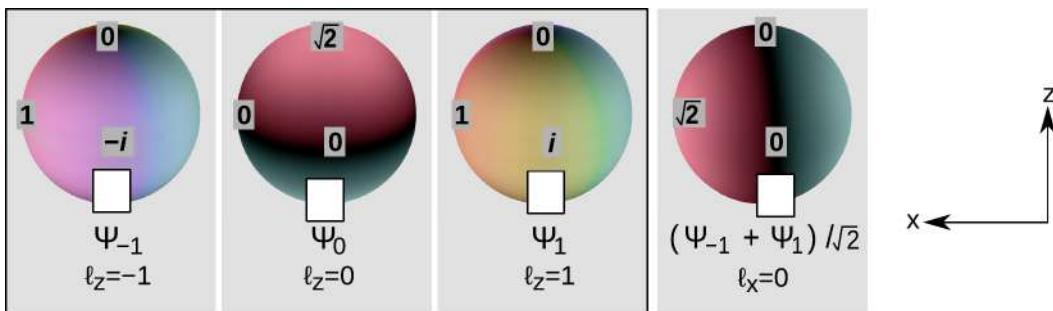
Determine the coefficients assigned to your group, and check your equation at the value of θ assigned to your group: $\theta = 0, \pi/2, \pi$, or $3\pi/2$.

4. An electron is initially in state Ψ_c , and Jane then measures its angular momentum. Discuss what happens to the electron's wavefunction and to Jane's, as in sec. 13.2.4, p. 887, on entanglement.

5. Discuss the probability interpretation and normalization.

Exercise 13C: Rotation around different axes

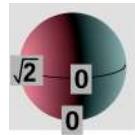
This exercise refers to the example at the beginning of section 13.4.3 on p. 925, which analyzes the figure below:



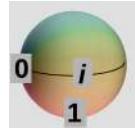
To simplify the writing:

- Ψ_{-1} means the state with $\ell_z = -1$, Ψ_0 has $\ell_z = 0$, etc.
- States with definite values of ℓ_x are notated as $\Psi_{\ell_x=0}$, etc.

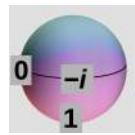
1. The wavefunctions are shown with values of the wavefunction written at the north pole and at two points on the equator. Fill in the south poles. Why do the results make sense physically for the $\ell_z = 1$ and -1 wavefunctions?
2. By rotating the pictures 90 degrees counterclockwise, we can make states of definite ℓ_x . We now want to express the states of definite ℓ_x in terms of the states of definite ℓ_z .



$$\Psi_{\ell_x=0} = \boxed{1/\sqrt{2}} \Psi_{-1} + \boxed{0} \Psi_0 + \boxed{1/\sqrt{2}} \Psi_1 \quad \text{Done on p. 925.}$$



$$\Psi_{\ell_x=1} = \boxed{} \Psi_{-1} + \boxed{} \Psi_0 + \boxed{} \Psi_1 \quad \text{Demonstrated by the instructor.}$$



$$\Psi_{\ell_x=-1} = \boxed{} \Psi_{-1} + \boxed{} \Psi_0 + \boxed{} \Psi_1 \quad \text{Done by the students.}$$

3. Fred takes a molecule known to have $\ell = 1$, and measures its ℓ_x . (This can be done by passing it through a magnetic field, as described in more detail in section 14.1.) If the molecule is not prepared in any particular orientation, then the result is random, and can be $\ell_x = -1$, 0, or 1. (The probabilities are all $1/3$, although this is not obvious.) Suppose he measures $\ell_x = 0$, so that *after* measurement, he is sure that the wavefunction is $\Psi_{\ell_x=0}$. (Fred may now be superimposed with other versions of himself who saw $\ell_x = -1$ or 1, but we stop keeping track of them now.)

Now suppose that Fred follows up with a second measurement, on the same molecule, but this time he orients the magnetic field so that he's measuring ℓ_z . What are the probabilities of the three possible results? Check normalization.

Exercise 13D: The Einstein-Podolsky-Rosen paradox

A nucleus having zero angular momentum undergoes symmetric fission into two fragments, each with $\ell = 1$. By conservation of momentum, they fly off back to back, and by conservation of angular momentum their angular momenta are also opposite. Let's say that except for this correlation, the two angular momentum vectors are randomly oriented.

1. Warm-up: Suppose Alice measures the ℓ_x of particle A, and Bob measures ℓ_x of fragment B. Make a table of the probabilities of the outcomes.

		particle B		
		$\ell_x = -1$	$\ell_x = 0$	$\ell_x = 1$
particle A	$\ell_x = -1$			
	$\ell_x = 0$			
	$\ell_x = 1$			

2. In a 1935 paper that ended up being one of the most frequently cited physics papers of all time, Einstein and his collaborators considered a scenario similar to the following. Suppose now that Alice measures ℓ_x , but Bob measures ℓ_z . It shouldn't matter who goes first, but let's say that Alice does. Using the results of exercise C, compute the probabilities in the row assigned to your group. Take into account the factor of $1/3$, because in this table, as in the first one, we're talking about the probability of a certain result for A *and* a certain result for B.

		Bob's probabilities		
		$\ell_z = -1$	$\ell_z = 0$	$\ell_z = 1$
1/3 of the time, Alice gets $\ell_x = -1 \implies$				
1/3 of the time, Alice gets $\ell_x = 0 \implies$				
1/3 of the time, Alice gets $\ell_x = 1 \implies$				
total probabilities for Bob				

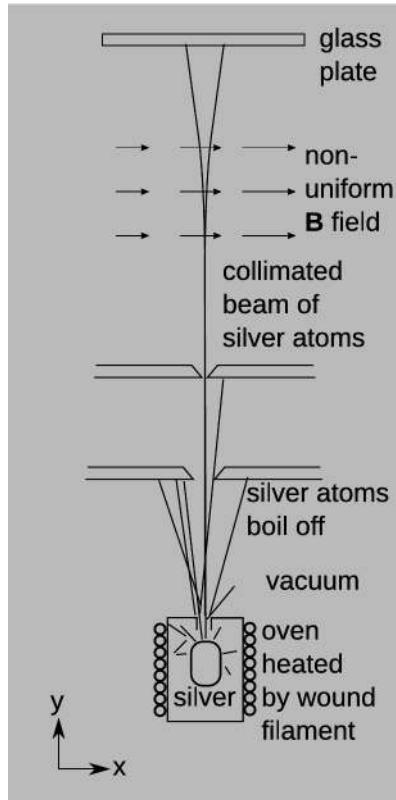
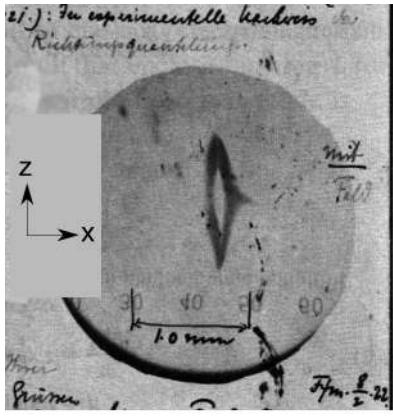
3. Can Alice send information to Bob by deciding whether or not to measure her particle's ℓ_x ?

Chapter 14

Additional Topics in Quantum Physics

14.1 The Stern-Gerlach experiment

In 1921, Otto Stern proposed an experiment about angular momentum, shown in figure a on p. 960, that his boss at the University of Frankfurt and many of his colleagues were certain wouldn't work. At this time, quantization of angular momentum had been proposed by Niels Bohr, but most physicists, if they had heard of it at all, thought of the idea as a philosophical metaphor or a mathematical trick that just happened to give correct results. World War I was over, hyperinflation was getting under way in Germany (a paper mark was worth a few percent of its prewar value), and the Nazi coup was still in the future, so that Stern, a Jew, had not yet been forced to flee to America. Because of the difficult economic situation, Stern and his colleague Walther Gerlach scraped up some of the funds to carry out the experiment from US banker Henry Goldman, cofounder of the investment house Goldman-Sachs.



a / Bottom: A schematic diagram of the Stern-Gerlach experiment. The z direction is out of the page. The entire apparatus is about 10 cm long. Top: A portion of Gerlach's celebratory 1922 postcard to Niels Bohr, with a photo showing the results. A coordinate system is superimposed. The orientation is flipped downward by 90 degrees compared to the schematic. The photo was taken through a microscope, and Gerlach drew the 1.0 mm scale on after the magnified photo had been printed.

The entire apparatus was sealed inside a vacuum chamber with the best vacuum obtainable at the time. A sample of silver was heated to 1000°C, evaporating it. The atoms leaving the oven encountered two narrow slits, so that what emerged was a beam with a width of only 0.03 mm, or about a third of the width of a human hair. The atoms then encountered a magnetic field. Because the atoms were electrically neutral, we would normally expect them to be unaffected by a magnetic field. But in the planetary model of the atom, we imagine the electrons as orbiting in circles like little current loops, which would give the atom a magnetic dipole moment \mathbf{m} . Even if we are sophisticated enough about quantum mechanics not to believe in the circular orbits, it is reasonable to imagine that such a dipole moment would exist. When a dipole encounters a *nonuniform* field, it experiences a force (example 7, p. 591). In this example, the forces in the x and z directions would be $F_x = \mathbf{m} \cdot (\partial\mathbf{B}/\partial x)$ and $F_z = \mathbf{m} \cdot (\partial\mathbf{B}/\partial z)$. (Because of Gauss's law for magnetism, these two derivatives are not independent — we have $\partial B_x/\partial x + \partial B_z/\partial z = 0$.) The rapidly varying magnetic field for this experiment was provided by a pair of specially shaped magnet poles, described in example 27, p. 745.

Because electrons have charge, we expect the motion of an electron to give it a magnetic dipole moment \mathbf{m} . But they also have mass, so for exactly the same reasons, we expect there to be some angular momentum \mathbf{L} as well. The analogy is in fact mathematically exact, as discussed in sec. 11.2.4, p. 697. Therefore this experiment with dipoles and magnetic fields is actually a probe of the behavior of angular momentum at the atomic level.

Luckily for Stern and Gerlach, who had no modern knowledge of atomic structure, the silver atoms that they chose to use do happen to have nonzero total \mathbf{L} , and therefore nonzero \mathbf{m} . The atoms come out of the oven with random orientations.

Classically, we would expect the following. Each atom has an energy $\mathbf{m} \cdot \mathbf{B}$ due to its interaction with the magnetic field, and this energy is conserved, so that the component m_x stays constant. However, there is a torque $\mathbf{m} \times \mathbf{B}$, and this causes the direction of the atom's angular momentum to precess, i.e., wobble like a top, with its angular momentum forming a cone centered on the x axis (example 25, p. 289). This precession is extremely fast, carrying out about 10^{10} wobbles per second, so that the atom precesses about 10^6 times while traveling the 3.5 cm length of the spectrometer. So even though the forces F_x and F_z are typically about the same size, the rapid precession causes F_z to average out to nearly zero, and only a deflection in the x direction is expected. Because the orientations of the atoms are random as they enter the magnetic field, they will have every possible value of L_x ranging from $-|\mathbf{L}|$ to $+|\mathbf{L}|$, and therefore we expect that when the magnetic field is turned on, the effect should be to smear out the image on the glass plate from a

vertical line to a somewhat wider oval. The atoms are dispersed from left to right along a certain scale of measurement according to their random value of L_x . The spectrometer is a device for determining L_x , a continuously varying number.

But that's all the classical theory. Quantum mechanically, L_x is quantized, so that only certain very specific values of F_x can occur. Although the discussion of precession above is really classical rather than quantum-mechanical, the result of F_z averaging out to zero turns out to be approximately right if the field is strong. Therefore we expect to see well separated vertical bands on the glass plate corresponding to the quantized values of L_x . This is approximately what is seen in figure a, although the field rapidly weakens outside the x - y plane, so we get the slightly more complicated pattern like a sideways lipstick kiss. Since we observe two values of L_x (the two "lips"), we conclude from these results that a silver atom has spin 1/2, so that L_x takes on the values $-\hbar/2$ and $+\hbar/2$. Although it took about five years for the experiment to be interpreted completely correctly, we now understand the Stern-Gerlach experiment to be not just a confirmation of the quantization of angular momentum along any given axis but also the first experimental evidence that the electron has an intrinsic spin of 1/2.

Discussion Questions

A Could the Stern-Gerlach experiment be carried out with a beam of electrons?

B A few weeks after the Stern-Gerlach experiment's results became public, Einstein and Ehrenfest carried out the following reasoning, which seemed to them to make the results inexplicable. Before a particular silver atom enters the magnetic field, its magnetic moment \mathbf{m} is randomly oriented. Once it enters the magnetic field, it has an energy $\mathbf{m} \cdot \mathbf{B}$. Unless there is a mechanism for the transfer of energy in or out of the atom, this energy can't change, and therefore the magnetic moment can only precess about the \mathbf{B} vector, but the angle between \mathbf{m} and \mathbf{B} must remain the same. Therefore the atom cannot align itself with the field. (They considered various mechanisms of energy loss, such as collisions and radiation, and concluded that all of them were too slow by orders of magnitude to have an effect during the atom's time of flight.) It seemed to them that as soon as the atom left the oven, it was somehow required to have anticipated the direction of the field and picked one of two orientations with respect to it. How can this paradox be resolved?

C Suppose we send a beam of oxygen molecules, with $L = \hbar$, through a Stern-Gerlach spectrometer, throwing away the emerging parts with $\ell_x = -1$ and +1 to make a beam of the pure $\ell_x = 0$ state. Now we let this beam pass through a second spectrometer that is identical but oriented along the z axis. Can we produce a beam in which every molecule has both $\ell_x = 0$ and $\ell_z = +1$? Hint: See the example in fig. d, p. 925.

14.2 Rotation and vibration

14.2.1 Types of excitations

Figure a shows the visible-light spectrum of the molecule N_2 . Because this particular chemical bond is unusually strong, the molecule does not break apart, even at the high temperature of a gas discharge tube, so we see the spectrum of the molecule, not of monoatomic nitrogen. This spectrum is more complex than the spectrum of the hydrogen atom, and that's not surprising, because the number of different states grows exponentially with the number of particles (here, 14 electrons plus 2 nuclei).

a / Visible spectrum of N_2 . Violet is on the left, red on the right.



What is more surprising is that there are some clear, simple patterns in this spectrum — patterns simpler than any that we would see in the spectrum of a monoatomic gas with the same number of particles. To start to understand this, we note that N_2 lacks the spherical symmetry of an individual atom, but it does have an axis of symmetry, b/1. These properties are also possessed by many nuclei, e.g., b/2. We now consider three different ways in which an excited energy state could occur in N_2 :

- Individual *particles* (electrons) can be raised to a higher energy level.
- The molecule can *vibrate* along its long axis, so that the nuclei (which have nearly all the inertia) move back and forth, elongating and compressing the system.
- The molecule can *rotate*.

14.2.2 Vibration

Particle excitations would produce the type of very complex, disorganized spectrum that we normally see in atoms that have many electrons, so that isn't what we're seeing in figure a. What about vibrations? For a classical harmonic oscillator, we know that the frequency of vibration is independent of the amplitude. If a classical oscillator contains electric charge, it will emit electromagnetic radiation at this frequency, smoothly and continuously draining itself of energy. As the energy is lost, the frequency stays the same.

b / 1. The molecule N_2 . 2. The nucleus ^{178}Hf .

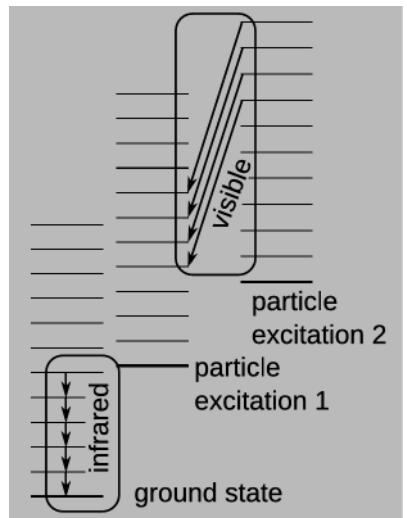
By the correspondence principle, we expect that when the quantum mechanical version of such a system is highly excited, it should emit a large number of photons of this frequency f , so that the discrete quantum jumps are undetectable and the radiation appears as a classical wave. We can thus infer that for a quantum vibrator, the excited states should show an *evenly spaced* ladder of energy levels.

Figure c shows how the series of red lines in figure a arises. For an excitation consisting only of vibrational motion, we expect based on the correspondence principle to see the evenly spaced ladder of states shown in a stack built above the ground state, with all of the photons having the same energy. These states and transitions do exist, but the light lies in the infrared spectrum and so is not seen in figure a. The red visible-light lines arise as shown in the second box, from states that involve both a certain particle excitation and some vibration. Because the spacing of the two ladders is slightly unequal, the red lines all have slightly different wavelengths.

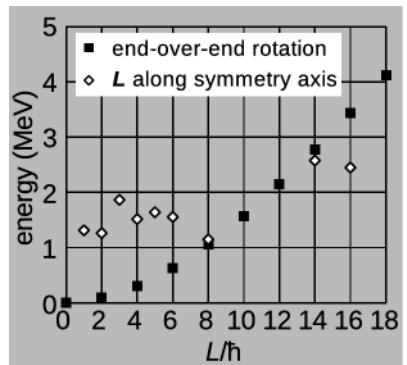
14.2.3 Rotation

What about rotation? An interesting thing happens here due to the structure of quantum mechanics. Quantum mechanics can describe motion only as a wave, with the value of the wave oscillating from one place to another. But this implies that according to quantum mechanics, no object can rotate about one of its axes of symmetry, for the rotated version of a state would then be the same state. This is why rotational excitations are never seen in individual atoms, or in nuclei that have spherical shapes. In examples like the ones in figure b, which have a single axis of symmetry, we can therefore have end-over-end rotation, but never rotation about the symmetry axis. Such end-over-end rotational states are observed in N_2 , for example, but because this involves large motions by the high-mass nuclei, the moment of inertia I is quite large, and therefore the rotational energies — classically, $K = L^2/2I$ — are very small, and infrared rather than visible photons are emitted. If rotation about the symmetry axis were possible, then the moment of inertia would be thousands of times smaller, because in such a rotation the nuclei would not move. The energies involved would be thousands of times higher, and the photons would lie approximately in the visible region of the spectrum. No such visible lines are actually observed.

Perhaps more vivid evidence for the nonexistence of rotation about a symmetry axis is shown in figure d. The states involving end-over-end rotation of the nucleus as a whole (“collective” rotation) are approximately a parabola on this graph, which is reasonable given the classical relation $K = L^2/2I$. But angular momentum cannot be generated along the symmetry axis through collective rotation. Instead, we see an irregular set of energy levels in which first one particle (for $L \leq 8\hbar$) and then two (14 and $16\hbar$) are excited.



c / Energy levels of the N_2 molecule.



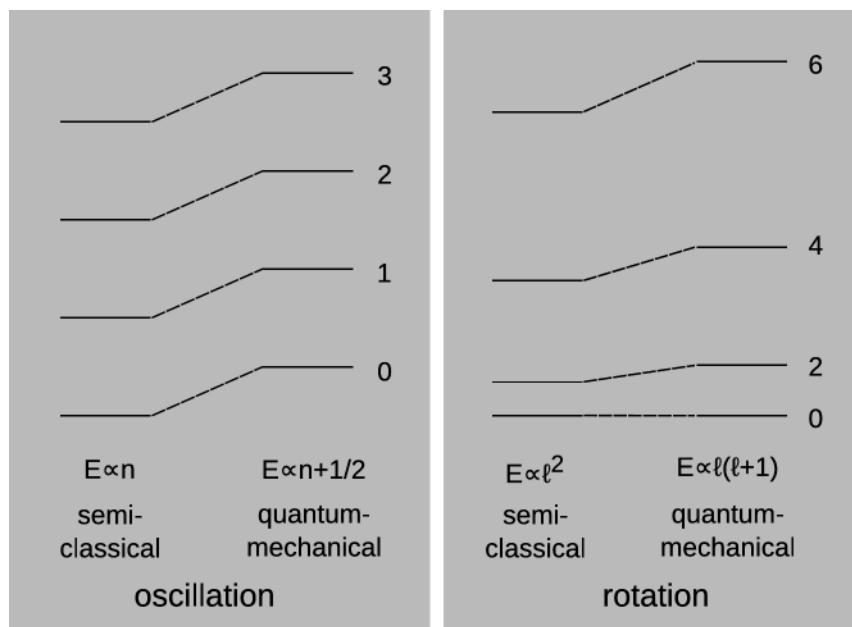
d / Excited states of the nucleus ^{178}Hf . Black squares represent states that are interpreted as end-over-end rotation, while white diamonds show particle excitations. For each angular momentum, the graph shows the lowest-energy state of each type, where known.

Note that only even multiples of \hbar are observed in collective rotation in figure d. This is because the nucleus's shape has an additional mirror symmetry, so that it is unaffected by a 180-degree rotation. This means that the wavefunction describing the collective rotation must oscillate twice as we pass through a full rotation.

14.2.4 Corrections to semiclassical energies

So far we've been using the correspondence principle to make educated guesses about quantum-mechanical expressions for the energies of vibrators and rotors. This style of reasoning is called semiclassical, because it combines ideas from classical and quantum physics. These expressions are guaranteed to be good approximations in the classical limit obtained when the quantum numbers are large, but figure e shows that the approximations can be poorer when the quantum numbers are small.

e / Quantum-mechanical corrections to the semiclassical results for the energy of a vibrator and a rotor. The rotational levels are shown for the case of a rotor with mirror symmetry, so that only even values of ℓ occur.



In the case of the n th excited state of a vibrator, the energy is $(n + 1/2)\hbar\omega$, where the $+1/2$ term represents a quantum correction to the semiclassical approximation. This shifts the entire ladder upward in energy by half a step. In particular, the energy of the ground state is not zero but rather $(1/2)\hbar\omega$. This can be verified quantitatively by calculating the energy for the solution to the Schrödinger discussed using the guess-and-check method in problem 23, p. 946. It is easy to see why the answer cannot be zero, for if it were, then the particle in the ground state would have zero kinetic energy and zero potential energy. To have zero kinetic energy, it would have to have a momentum of exactly zero, so $\Delta p = 0$, but to have zero potential energy it would also have to sit still at exactly the equi-

librium position, so $\Delta x = 0$. But this would violate the Heisenberg uncertainty principle and so is impossible.

The inevitable motion that is present even in the ground state is known as zero-point motion, and its energy is the zero-point energy. Relativity tells us that $E = mc^2$, so the zero-point energy of particles is equivalent to a certain amount of mass. In fact, nearly all the mass of ordinary matter arises from the zero-point energy of the quarks inside the neutrons and protons. Another interesting application is to spontaneous nuclear fission, which is the basis for nuclear energy, providing the kick-off for a chain reaction. Spontaneous fission requires that a nucleus become more and more elongated until it breaks apart into two pieces. The very elongated shapes have a high potential energy, so that spontaneous fission requires quantum-mechanical tunneling. If it were not for the zero-point vibrational energy associated with this motion, the tunneling probability for uranium and plutonium isotopes would be extremely small. These isotopes would decay only by alpha emission, and nuclear reactors and bombs would not work.

Chemistry of deuterium

example 1

Chemistry is an electrical interaction, and neutrons have no charge, so to a first approximation we expect the number of neutrons in a nucleus to have no effect on chemical properties. That is, all isotopes of an element are typically expected to have the same chemistry. But there are some unusual cases where different isotopes can have rather different chemical behavior, an example being the differences between hydrogen-1 (ordinary hydrogen) and hydrogen-2, also known as deuterium, ${}^2\text{H}$ or just D for short. The masses of these isotopes differ by a factor of 2, which is unusually large for two isotopes of the same element.

As an example, we consider the bond between a carbon atom and a hydrogen atom, which is important in organic and biochemistry. Classically, we can imagine this as a big mass and a small mass, joined by a spring. Infrared spectroscopy of the molecule CH, with ordinary hydrogen, shows that when the $n = 1$ vibrational state emits a photon and decays to the $n = 0$ ground state, the photon has $\hbar\omega_{\text{photon}} = 0.3521 \text{ eV}$. By conservation of energy, this equals the difference between the vibrational states of the molecule $E_1 - E_0 = (1 + 1/2)\hbar\omega - (0 + 1/2)\hbar\omega = \hbar\omega$, i.e., the frequency of the photon is the same as the frequency of the molecular vibration (as we would expect classically, since charges oscillating at a certain frequency will produce radiation at that frequency).

Now classically, the frequency of a simple harmonic oscillator is $\omega = \sqrt{k/m}$, where k is the spring constant and m is the mass. The “spring constant” k here is the stiffness of the bond, which arises from electrical interactions, and is therefore identical for ordinary CH and for the CD molecule formed with deuterium. Be-

cause the mass of the H or D is an order of magnitude smaller than that of the carbon, it's a pretty good approximation to say that the carbon stands still while the H or D vibrates, and therefore we can approximate the mass m as being just the mass of the H or D. We therefore expect ω to differ by a factor of about $\sqrt{2}$ between CH and CD, with the latter frequency being the lower one. This means that the zero-point energy $(1/2)\hbar\omega$ differs by a factor of $\sqrt{2}$. The difference works out to be about 0.05 eV, the energy being smaller for CD. We therefore expect that a CD bond will be more stable than a CH bond by this amount of energy. This is a substantial amount, so in organic chemistry, we expect that there will be a nonnegligible difference in behavior. Deuterium atoms will tend to displace hydrogen atoms.

In addition to this effect, there are also significant differences in the behavior of hydrogen and deuterium in their ability to tunnel through a potential barrier, which can be an important effect when protons are transferred between molecules.

For these reasons, there can be serious consequences if a living organism is given deuterium-containing water ("heavy water") to drink. If the percentage of deuterium becomes a significant fraction of all hydrogen isotopes in the body of a multicellular organism, its metabolism is disrupted enough to kill it. This is surprising, since we ordinarily expect no chemical effects from substituting one isotope for another.

f / Each panel of the figure shows a standing wave on a sphere, with the convention that gray is zero, white is a positive real number, and black is a negative real number. (These could instead have been drawn as traveling waves, but then we would have needed to represent complex numbers using color, as in figure c on p. 924.) Only 2 is a solution of the Schrödinger equation.

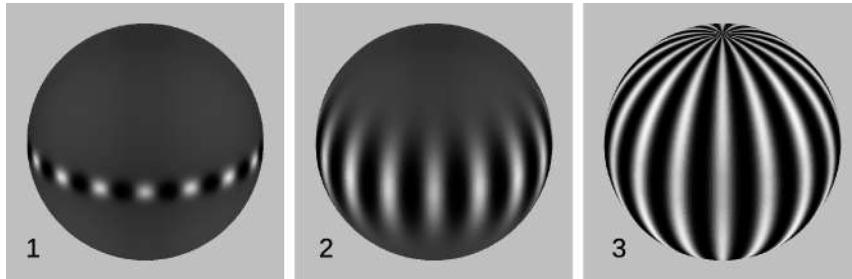


Figure f shows visually the reason for the correction of ℓ^2 to $\ell(\ell+1)$. Each of these standing waves has $|\ell_z| = 16$, where z is the vertical axis. But only f/2 is a solution of the Schrödinger equation for a state of definite ℓ . To be a solution of the Schrödinger equation, such a wave must have the same kinetic energy everywhere. Each of these three has the same kinetic energy associated with its wavelength in the "east-west," or azimuthal, direction. Wave f/1 is not a solution, because near the equator, it has an extremely short wavelength in the "north-south," or longitudinal, direction, and this gives it a greater kinetic energy near the equator than elsewhere. The opposite problem occurs in f/3, where the wave is constant in the longitudinal direction; at the poles, the wavefunction varies in-

finitely rapidly, and therefore the kinetic energy blows up to infinity there. The only valid solution is $f/2$, which has a Goldilocks-style just-right wavelength in the longitudinal direction. The kinetic energy associated with this wavelength is the difference between the semiclassical ℓ^2 and the correct quantum mechanical $\ell(\ell + 1)$.

A different example that is particularly easy to reason about is the wavefunction Ψ_{10} shown in figure d on p. 925, for $\ell = 1$ and $\ell_z = 0$. (The odd value of ℓ is possible for a rotor that doesn't have mirror symmetry, e.g., the carbon monoxide molecule CO.) The ratio of the correct quantum mechanical energy to the semiclassical one is $\ell(\ell + 1)/\ell^2 = 2$, and the factor of two makes sense because at the poles, the wave has equal contributions to its kinetic energy due to oscillations in the two perpendicular directions that occur in the Laplacian $\partial^2/\partial x^2 + \partial^2/\partial y^2$.

Discussion Question

A The correction of the semiclassical proportionality for the energy of a rotor from ℓ^2 to $\ell(\ell + 1)$ is effectively the addition of a correction equal to ℓ . What if someone tells you that there is an additional correction term that depends only on ℓ_z (for a fixed ℓ)? Is this plausible?

B Can the correction $\ell^2 \rightarrow \ell(\ell + 1)$ be tested experimentally by measuring the energy of a spinning steel ring in the laboratory? Can the correction $n \rightarrow n + 1/2$ be tested using a cart on an air track that vibrates back and forth between two springs?

14.3 ★ A tiny bit of linear algebra

This optional section is a self-contained presentation of a very small amount of linear algebra. None of the later physics requires this material, but reading it may be helpful as review for the reader who has already had an entire linear algebra course, or to help make connections for the one who is taking such a course concurrently or will take it in the future.

A *vector space* is a set of objects, which we refer to as vectors, along with operations of addition and scalar multiplication defined on the vectors. The scalars may be the real numbers or the complex numbers. We require that the addition and scalar multiplication operations have the properties that addition is commutative ($\mathbf{u} + \mathbf{v} = \mathbf{v} + \mathbf{u}$), that we have an additive identity 0 and additive inverses ($\mathbf{v} + (-\mathbf{v}) = 0$), and that both operations are associative and distributive in the ways that we would expect from the notation. The prototypical example of a vector space is vectors in three-dimensional space, with the scalars being the real numbers.

The vector space of polynomials

example 2

Consider the set of all polynomials. If we define addition of polynomials and multiplication of a polynomial by a real number in the obvious ways, then these functions are a vector space. Note that

there is no well-defined division operation, since dividing a polynomial by a polynomial typically does not give a polynomial.

In quantum mechanics, we are interested in the vector space of wavefunctions, with the scalars being the complex numbers.

A set of vectors is said to be *linearly independent* if it is not possible to form the zero vector as a linear combination of them. For vectors in three-dimensional space, a set of three vectors is not linearly independent if they lie in the same plane. The set of polynomials $\{1, x\}$ is linearly independent, but the set $\{P, Q, R\}$, where $P = 1$, $Q = 1 - x$, and $R = 1 + x$, is not, because $-2P + Q + R = 0$.

A *basis* for a vector space is a linearly independent set of vectors, called basis vectors, such that any vector can be formed as a linear combination of basis vectors. The standard basis for vectors in two-dimensional space is $\{\hat{x}, \hat{y}\}$, while a possible basis for the polynomials is the infinite set $\{1, x, x^2, x^3, \dots\}$. A basis exists for any vector space, and in fact there are normally many different bases to choose from, with none being preferred. In the plane, for example, we can choose to rotate the standard $\{\hat{x}, \hat{y}\}$ basis by any angle we like. Every basis for a given vector space has the same number of elements, and this number is called the *dimension* of the vector space. The plane is a two-dimensional vector space. The polynomials are an infinite-dimensional vector space.

A *linear operator* is a function \mathcal{O} that takes a vector as an input and gives a vector as an output, with the properties $\mathcal{O}(\mathbf{u} + \mathbf{v}) = \mathcal{O}(\mathbf{u}) + \mathcal{O}(\mathbf{v})$ and $\mathcal{O}(\alpha\mathbf{u}) = \alpha\mathcal{O}(\mathbf{u})$. A rotation in the plane is a linear operator.

Differentiation as a linear operator

example 3

Consider the set of all differentiable functions, taken as a vector space over either the real numbers or the complex numbers. Then the derivative is a linear operator, as is the second derivative. The kinetic energy term in the Schrödinger equation is built out of second derivatives, so it is a linear operator.

For vectors in three-dimensional space, we have a dot product, which is a function that takes two vectors as inputs and gives a scalar as its output. A vector space may or may not come equipped with such an operation. If it does, we call the operation an *inner product*. The inner product on wavefunctions is introduced in section 14.6.2, p. 984. In quantum mechanics, the inner product is a basic tool used to define probabilities, and for example normalization becomes the requirement that a wavefunction have an inner product with itself that equals 1. That is, a normalized wavefunction is a kind of unit vector.

When a vector space is finite-dimensional and a basis has been chosen, then if we wish we can represent vectors in column vector notation. For example, in the space of first-order polynomials with

the basis $\{1, x\}$, the polynomial $3 + 5x$ can be represented by $(\begin{smallmatrix} 3 \\ 5 \end{smallmatrix})$. Linear operators can similarly be represented by matrices, but we will seldom find this possible or useful in this book. For example, we can't represent the derivative as a matrix, because the vector space is infinite-dimensional.

14.4 The underlying structure of quantum mechanics, part 1

So far we have been building up the structure of quantum mechanics by casually laying one brick on top of another, but at this point it will be advantageous to pause and consider the broader blueprint.

14.4.1 The time-dependent Schrödinger equation

For simplicity, our discussion of the Schrödinger equation in section 13.3.6, p. 907, was limited to standing waves, allowing us to avoid explicitly discussing how the wavefunction changed with passing time. Let's consider the generalization to the full time-dependent case.

Classically, suppose I show you a picture of a baseball next to a tree, and I ask you how long it will take to hit the ground. You can't tell, because you also need information about the ball's initial velocity. That is, the future time-evolution of the system $x(t)$ depends not just on the initial position $x(0)$ but also on its initial time derivative $x'(0)$.

But if I show you a uranium atom in its lowest energy state, you don't need to know any other information to predict everything that can be predicted about its future decay. Whereas the baseball could be thrown downward in order to make it reach the ground more quickly, nobody knows of any way to prepare the uranium nucleus in such a way that it is any more likely to decay sooner. Knowing the initial wavefunction $\Psi(0)$ to be that of the ground state lets us say as much as can be said about the future time-evolution $\Psi(t)$, and it's neither necessary nor helpful to know the time derivative $\Psi'(t)$.

This is an example of a more general idea about the interpretation of quantum mechanics, which is that the wavefunction is a complete description of any system. There isn't more information that can be known about the system. This principle seems to be widely agreed upon by physicists, but doesn't seem to have a standard name. (The phase and normalization of the wavefunction are not considered to give any information, since the phase is unobservable, and the normalization can be standardized so that the total probability is 1. See the sidebar for more detail.)

Unobservability of phase and normalization

When we say that phase and normalization don't count as knowledge of a system, we're saying something very mathematically specific: that Ψ and $c\Psi$ represent the same state, where $c \neq 0$ is a complex number; the magnitude of c would only affect the normalization, and its argument would only affect the phase. We do not mean, for example, that wavefunctions like $\sin x$ and $\cos x$ are indistinguishable. The sine and cosine give different probability distributions, so they are distinguishable. For example, the $\sin x$ wavefunction gives zero probability of detection at $x = 0$. See also problem 17, p. 1015 and example 9, p. 982.

Linear algebra application

Wavefunctions can be described by vectors in a vector space (p. 967). A state is a one-dimensional subspace of the vector space, i.e., the set of all wavefunctions of the form $c\Psi$ for some fixed Ψ .

Wavefunction fundamentalism

All knowable information about a system is encoded in its wavefunction (ignoring phase and normalization).

An example of an idea that would violate this principle is the pilot wave theory proposed by de Broglie around 1927, and improved by Bohm in the 1950's. In this theory, an electron-particle is a separate object from an electron-wave, with the particle surfing the wave along a deterministic trajectory.

Another example that shows the contrast with the classical description is that if I show you a snapshot of a wave on a string, you can't tell which direction it's going — as with the baseball, you need to know its initial velocity in addition. But if I show you a snapshot of a quantum-mechanical traveling wave, you *can* tell which direction it's going, because of the complex phase, as shown in figures u/2 and u/3 on page 916. Note that this mechanism wouldn't work if wavefunctions were always real numbers, so wavefunction fundamentalism implies complex wavefunctions.

Given the wavefunction at some initial time, we can predict its evolution into the future by making use of the principle that $E = hf$. Suppose for example that we have a sinusoidal plane wave traveling to the right. Then we expect the value of the wavefunction at a particular point in space to rotate clockwise about the origin in the complex plane at the appropriate frequency f , showing a time dependence $e^{-i\omega t}$ (where, as usual, $\omega = 2\pi f$). Thus the time derivative of the wavefunction is $\Psi' = -i\omega\Psi = -i(E/\hbar)\Psi$, so that $E\Psi = i\hbar\Psi'$. Then to generalize the time-independent Schrödinger equation to its time-dependent version, the most obvious thing to try is simply to substitute $i\hbar\partial\Psi/\partial t$ for $E\Psi$, which gives

$$i\hbar\frac{\partial\Psi}{\partial t} = -\frac{\hbar^2}{2m}\nabla^2\Psi + U\Psi.$$

(In section 14.6.4, p. 991, we will generalize this to cases where the wavefunction is not expressed in terms of the spatial coordinates x , y , and z .) Unlike Newton's laws of motion, which refer to a second derivative with respect to time, the Schrödinger equation involves only a first time derivative. This is why we don't need initial data on $\partial\Psi/\partial t$, but only Ψ : if we know Ψ , then the right-hand side of the Schrödinger equation is what *gives* us $\partial\Psi/\partial t$. But the Schrödinger equation has some other properties that match up with those of Newton's laws.

A plane wave

example 4

Consider a free particle of mass m in one dimension, with the wavefunction

$$\Psi = e^{i(kx-\omega t)},$$

where $k = 2\pi/\lambda = p/\hbar$ is called the wavenumber. If k and ω are both positive, we can tell that the particle is moving to the right,

because the signs inside the exponential are such that x could increase as t increases while keeping the phase the same. This would happen for $k\Delta x - \omega\Delta t = 0$, or $v = \omega/k$, which is the phase velocity (not the same as the group velocity, sec. 13.3.2, p. 896).

Suppose that the particle is in free space, so that U is constant, and for convenience take $U = 0$. Application of the Schrödinger equation, $i\hbar\partial\Psi/\partial t = -(\hbar^2/2m)\partial^2\Psi/\partial x^2$, gives $\hbar\omega e^{(\dots)} = \frac{\hbar^2 k^2}{2m} e^{(\dots)}$, and if this is to hold true for all values of x and t , then we must have $\hbar\omega = \frac{\hbar^2 k^2}{2m}$, which is simply an expression of the Newtonian relation $K = p^2/2m$, since $k\lambda = 2\pi$ and $p = h/\lambda$. Flipping the sign of k results in an equally valid solution, and a negative k is how we would represent a wave traveling to the left.

We have two solutions to the Schrödinger equation corresponding to the two signs of k , and because the Schrödinger equation is linear, it follows that we can make a more general solution of the form

$$Ae^{i(kx-\omega t)} + Be^{i(-kx-\omega t)},$$

where A and B are any two complex numbers. (We could also try to elaborate on this theme by allowing for an arbitrary phase angle δ inside the complex exponentials, e.g., changing the argument of the first exponential to $i(kx - \omega t + \delta)$. However, this would be equivalent to changing A to $Ae^{i\delta}$, which is just a change in A 's phase angle, not a new solution.)

Dispersion of a wave packet

example 5

An annoying feature of example 4 is that the wavefunction cannot be normalized because it extends in all directions to infinity. This type of infinite plane wave is at best an idealization of the wavefunction for a realistic example such as an electron launched by a cathode ray tube, or an alpha particle emitted by a nucleus. As a more realistic example, we might try something like a wave packet, such as a pulse with a certain shape, traveling in a certain direction. This works for waves on a string or for electromagnetic waves: the pulse or packet simply glides along while rigidly maintaining its shape. To investigate this idea using the time-dependent Schrödinger equation, it will be convenient to adopt the frame of reference in which the particle is at rest. In this frame, we would expect the wave to look frozen, just as an ocean wave looks frozen in place to a surfer who is riding it. It must therefore be of the form

$$\Psi = e^{-i\omega t} f(x),$$

where f is some function specifying the shape of the wave packet. But this Ψ is not a solution to the Schrödinger equation. On the left-hand side of the Schrödinger equation, evaluating the time derivative gives $\hbar\omega\Psi$, which is just the original wavefunction multiplied by a constant. If we are to satisfy the Schrödinger equation, then the right-hand side, which is the second derivative

with respect to x , must also be equal to the original wavefunction multiplied by a constant. But the only functions for which $(d^2 / dx^2)(\dots) = (\text{constant})(\dots)$ are exponentials and sine waves. An exponential shape obviously isn't a physical realization of a wave packet. A sine wave works, but it just describes an infinite plane wave like the one in example 4, not a wave packet that can be localized and normalized.

The underlying reason for this result is that the Schrödinger equation is dispersive: waves with different wavelengths travel at different speeds (because they correspond to different momenta). Suppose a pulse has the shape $f(x)$ at $t = 0$. Since a pulse is not a sine wave, it doesn't have a single well-defined wavelength, and therefore it doesn't have a definite momentum or velocity. In fact, the spread in momentum must be at least a certain size due to the Heisenberg uncertainty principle $\Delta p \Delta x \gtrsim h$. This causes the pulse to spread out over time.

This leads to a strange thought experiment. Suppose that a uranium atom in the Andromeda galaxy emits an alpha particle, which then travels thousands of light years and eventually flies past the earth. Its wave packet may initially have been as narrow as the diameter of an atomic nucleus, $\sim 10^{-15}$ m, but by the time it arrives perhaps it is the size of an aircraft carrier. Will an observer see a gigantic alpha particle flying by? No, because observing it constitutes a measurement of its position, and by the probability interpretation of the wavefunction this measurement simply has a certain probability of giving a result that is anywhere within some region the size of an aircraft carrier.

14.4.2 Unitarity

The Schrödinger equation is completely deterministic, so that if we know Ψ initially, we can always predict it in the future. We can also “predict” backward in time, so that the system’s history can always be recovered from knowledge of its present state. Thus there is never any loss of information over time. Furthermore, it can be shown that probability is always conserved, in the sense that if the wavefunction is initially normalized, it will also be normalized at all later times.

Linear algebra application

Time evolution is represented by a linear operator (p. 968). Unitarity is an additional requirement for this linear operator.

Unitary evolution of the wavefunction

The wavefunction evolves over time, according to the Schrödinger equation, in a deterministic and *unitary* manner, meaning that probability is conserved and information is never lost.

(Unitarity is defined more rigorously on p. 987.)

Since we think of quantum mechanics as being all about randomness, this determinism may seem surprising. But determinism in the time-evolution of the wavefunction isn’t the same as deter-

minism in the results of experiments as perceived and recorded by a human brain. Suppose that you prepare a uranium atom in its ground state, then wait one half-life and observe whether or not it has decayed, as in the thought experiment of Schrödinger's cat (p. 887). There is no uncertainty or randomness about the wavefunction of the whole system (atom plus you) at the end. We know for sure what it looks like. It consists of an equal superposition of two states, one in which the atom has decayed and your brain has observed that fact, and one in which the atom has not yet decayed and that fact is instead recorded in your brain.

To get more of a feeling for what is meant by unitarity, it may be helpful to consider some examples of how it could be violated. One is the mythical “collapse” of the wavefunction in naive interpretations of the Copenhagen approximation (p. 889). Another example of nonunitarity is given in example 15 on p. 992.

A more exotic example is the disappearance of matter into a black hole. If I throw my secret teenage diary into a black hole, then it contributes a little bit to the black hole's mass, but the embarrassing information on the pages is lost forever. This loss of information seems to imply nonunitarity. This is one of several arguments suggesting that quantum mechanics cannot fully handle the gravitational force. Thus although physicists currently seem to possess a completely successful theory of gravity (Einstein's theory of general relativity) and a completely successful theory of the microscopic world (quantum mechanics), the two theories are irreconcilable, and we can only make educated guesses, for example, about the behavior of a hypothetical microscopic black hole.

14.5 Methods for solving the Schrödinger equation

14.5.1 Cut-and-paste solutions

Quite a few of the interesting phenomena of quantum mechanics can be demonstrated by finding solutions to the one-dimensional Schrödinger equation using the following “cut and paste” method. We break up the x axis into pieces, where the potential $U(x)$ does different things, and such that we already know the solutions of the Schrödinger equation for each piece. We then splice together the different parts of the solution, requiring that no discontinuities occur in the wavefunction Ψ or its derivative $\partial\Psi/\partial x$. (If the momentum and kinetic energy are to be finite, and U is finite, then we need all derivatives up to the second to be defined.)

Partial reflection at a step

The simplest example of this kind is a potential step,

$$U(x) = \begin{cases} U_1, & x < 0 \\ U_2, & x > 0, \end{cases}$$

where U_1 and U_2 are constants, and the energy of the particle is such that both sides are classically allowed. We have discussed this example 18 on p. 909, where we cheated by drawing real-valued wavefunctions, and simply assumed that we could still use our previous results for classical wave reflection (p. 381). It is not actually obvious that we should be able to get away with recycling that result, both because our quantum-mechanical wavefunctions are complex and because the Schrödinger equation is dispersive, so we can no longer assume, as we did there, that a wave packet simply glides along rigidly (example 5, p. 971).

To sidestep the problem of dispersion, we will carry out our analysis using an infinitely long wave-train with a definite wavelength. Let the incident wave have unit amplitude and travel to the right,

$$\Psi_I = e^{i(kx - \omega t)} \quad (x < 0),$$

as in example 4, p. 970. Recall that the wavenumber k is basically just momentum, $p = \hbar k$.

For the reflected and transmitted parts of the wave, we take

$$\Psi_R = R e^{i(-kx - \omega t)} \quad (x < 0),$$

and

$$\Psi_T = T e^{i(k'x - \omega t)} \quad (x > 0),$$

where the reflected and transmitted amplitudes R and T are unknown, and our goal is to find them. The different sign inside the exponential for Ψ_R corresponds to the opposite direction of motion at the same speed v , while in the expression for Ψ_T we have motion to the right, but with a different momentum $p' = \hbar k'$ as required by conservation of energy.

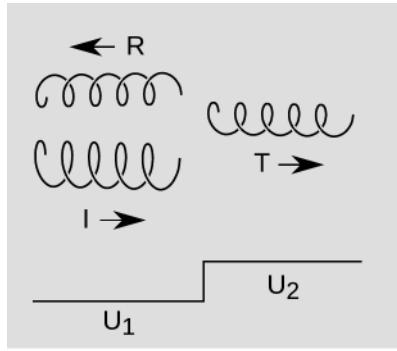
Demanding continuity of Ψ gives

$$1 + R = T.$$

The derivatives are $\partial\Psi_I/\partial x = ik\Psi_I$, $\partial\Psi_R/\partial x = -ik\Psi_R$, and $\partial\Psi_T/\partial x = ik'\Psi_T$, and evaluating these at $x = 0$, $t = 0$ gives ik , $-ikR$, and $ik'T$. If the derivative is to be the same for $x \rightarrow 0^-$ and for $x \rightarrow 0^+$, we need to have $ik - Rik = iTk'$, or

$$1 - R = \frac{k'}{k}T,$$

But these two equations are exactly the same as the ones found on p. 381 for a classical, nondispersive wave, the only difference being



a / An incident wave is partially reflected and partially transmitted at a step in the potential U . The complex wavefunctions are represented using a complex plane perpendicular to the direction of propagation, so that they look like corkscrews. The incident and reflected wavefunctions actually superposed, but are drawn as separate entities and offset for purposes of visualization.

the replacement of v/v' with k'/k . To keep the writing simple, let $\alpha = k'/k$. With this replacement, the solutions are the same as before, $R = (1 - \alpha)/(1 + \alpha)$ and $T = 2/(1 + \alpha)$. For a particle of energy E , we can find the momentum ratio α using conservation of energy, $\alpha = \sqrt{(E - U_2)/(E - U_1)}$. There is partial reflection not just in the case of a sudden rise in the potential, but also at a sudden drop ($U_2 < U_1$), which is surprising and seems to violate the correspondence principle, but actually does not, as discussed in example 18 on p. 909.

One of our principles of quantum mechanics is unitarity (p. 972), which says, in part, that probability is conserved. Normally we would interpret this to mean that a wavefunction stays normalized if it was originally normalized. In this example, the wavefunctions are not normalizable, but we still expect the fluxes of particles balance out. We have

$$\begin{aligned} \text{flux} &= (\text{probability density})(\text{group velocity}) \\ &= \Psi^2 \cdot \frac{p}{m} \\ &= \frac{\hbar}{m} k \Psi^2, \end{aligned}$$

so that if we want the total incident flux to equal the total outgoing flux, we need

$$k = kR^2 + k'T^2,$$

which is straightforward to verify.

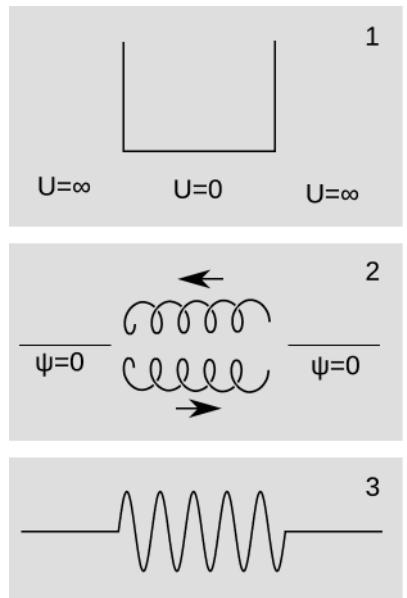
Infinite potential well

In sec. 13.3.3, p. 899, we analyzed the one-dimensional particle in a box. There was nothing wrong with those results, but it is of interest to see how they fit into the framework of the time-dependent Schrödinger equation. If we want the walls of the box to be completely impenetrable, then we should describe it using a potential such as

$$U(x) = \begin{cases} \infty, & x < 0 \\ 0, & 0 < x < L, \\ \infty, & x > L, \end{cases}$$

shown in figure b/1. Because the potential is infinite outside the box, we expect that there is no tunneling, and zero probability of finding the particle outside.

In general when we do the cut-and-paste technique, we expect both the wavefunction and its first derivative to be continuous where the pieces are joined together. But because we have already solved this problem by more elementary methods, we know that there will be kinks in the wavefunction at the walls of the box, $x = 0$ and L . The kink is a point where the second derivative $\partial^2\Psi/\partial x^2$ is undefined, and it's undefined because it's infinite. The second derivative



b / A particle in a box.

is essentially the kinetic energy operator, and normally it would not be possible to have the kinetic energy be $\pm\infty$. But in this problem, it *is* reasonable to have a kinetic energy of $-\infty$, because the potential energy is $+\infty$.

Within the box, for a fixed energy $E = \hbar\omega$, the possible wavefunctions will be those of a free particle, which we have already found. There are two possibilities, of the form

$$\begin{aligned}\Psi_1 &= e^{i(kx-\omega t)} \\ \Psi_2 &= e^{i(-kx-\omega t)},\end{aligned}$$

figure b/2, where k is a positive real number satisfying $k = p/\hbar = \sqrt{2mE}/\hbar$. Ψ_1 is a wave traveling to the right, and Ψ_2 is a wave traveling to the left. The most general solution will be a superposition of these,

$$\Psi = A\Psi_1 + B\Psi_2.$$

Because the wavefunction has to be continuous at $x = 0$, where $\Psi_1 = \Psi_2$, we must have $A + B = 0$. Since $e^{iz} - e^{-iz} = 2\sin z$, we end up with

$$\Psi = 2A \sin kxe^{-i\omega t}.$$

Throwing out the time-dependent phase, we get the sinusoidal solutions to the time-independent Schrödinger equation that we have already found, e.g., figure b/3. Imposing the additional constraint that Ψ be continuous at $x = L$, we get the condition $kL = n\pi$, where n is an integer, and this makes the energies quantized, as we found before.

14.5.2 Separability

When we first generalized the Schrödinger equations from one dimension to two and three dimensions, a trick for finding solutions was to take solutions to the one-dimensional equation and multiply them. For example, we knew that in the case of a constant potential (a free particle), the one-dimensional time-independent equation had solutions of the form $\sin ax$ and e^{ax} . We then saw in problem 37, p. 949, that $e^{by} \sin ax$ was a solution to the two-dimensional equation. This is because the two-dimensional time-independent Schrödinger equation for a free particle, which has the form

$$\nabla^2\Psi = c\Psi,$$

has a property called *separability*. What this means is that if functions X and Y are both solutions of the one-dimensional version of the equation, then $\Psi(x, y) = X(x)Y(y)$ is a solution of the two-dimensional one. To see this, we calculate

$$\begin{aligned}\nabla^2\Psi &= \nabla^2[X(x)Y(y)] \\ &= \left(\frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} \right) [X(x)Y(y)] \\ &= Y(y)X''(x) + X(x)Y''(y).\end{aligned}$$

We're looking for functions X and Y such that this is a solution to the two-dimensional equation, so that

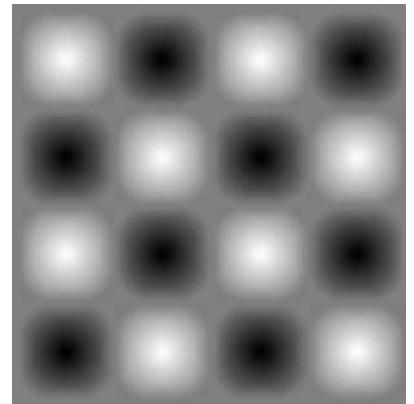
$$Y(y)X''(x) + X(x)Y''(y) = cX(x)Y(y).$$

Dividing both sides by $X(x)Y(y)$ simplifies this equation to

$$\frac{X''(x)}{X(x)} + \frac{Y''(y)}{Y(y)} = c.$$

But if X and Y are solutions of the one-dimensional equation, then both terms on the left are constants, so we have a valid solution to the two-dimensional equation.

As an example, we know that $\sin kx$ is a solution to the one-dimensional Schrödinger equation, so the function $\sin kx \sin ky$ is also a solution. The result, shown in figure c, can be chopped off and made into a solution of the two-dimensional particle in a box. Solutions similar to this one are found in real-life examples such as microwave photons in a microwave oven. For more about separability, and how it compares with entanglement, see sec. 14.11, p. 1007.



c / A solution to the Schrödinger equation found by separability. Positive values are shown as light colors, negative ones as dark colors.

14.6 The underlying structure of quantum mechanics, part 2

14.6.1 Observables

By the time my first-year mechanics students have been in class for a week, they know how to answer when I ask them the velocity of the tape dispenser at the front of the classroom: “We don’t know, it depends on your frame of reference.” The *absolute* velocity of an object is a meaningless concept, part of the mythical dungeons-and-dragons cosmology of Aristotelian physics. Quantum mechanics is as great a break from Newton as Newton was from Aristotle, and similar care is required in redefining what concepts are *observables* — meaningful things to talk about measuring.

Classically, we describe the state of the system as a point in phase space (sec. 5.4.2, p. 328) — which is just a fancy way of saying that we specify all the positions and momenta of its particles — and an observable is defined as a function that takes that point as an input and produces a real number as an output. (By the way, the word “phase” in “phase space” doesn’t refer directly to the phase of a wave, which we’ll also be discussing below.) For example, kinetic energy is a classical observable, and $K(\textcircled{O}) = 0$, where the picture represents a tennis ball at rest. For a moving tennis ball with one unit of energy, $K(\textcircled{E}\textcircled{O}) = 1$. For a vibrating violin string, we could have $U(\textcircled{W}) = 1$, and $U(\textcircled{V}) = 4$ (where doubling the amplitude gives four times the energy).

Quantum-mechanically, the Heisenberg uncertainty principle tells us that we can't independently dial in the desired values of a particle's position and momentum. They aren't two variables that are independent of one another. Therefore we don't have a phase space, so an observable has to be represented by a function whose input is a wavefunction. Furthermore, we expect that:

- The output shouldn't depend on the phase¹ of the wavefunction.
- The output shouldn't depend on amplitude (because a different amplitude might just mean an incorrectly normalized state).
- The output should be well defined when we superpose any two states.

These requirements are hard to reconcile with the idea that the output of the observable is just a real number representing the result of the measurement. We could decree that the input wavefunction is just required to be have the standard normalization, but there's no obvious way to define a standardization of phase. And suppose we have a particle in a one-dimensional box, with the two lowest energies being $E(\curvearrowleft) = 1$ and $E(\curvearrowright) = 4$. Then what should we define for the superposition $E(\curvearrowleft + \curvearrowright)$? We could define it to be the average, 2.5, but that isn't even a possible value of the measurement; in reality, the result of the measurement would be either 1 or 4, with equal probability.

For a clue as to a better way to proceed, note the structure of the time-independent Schrödinger equation for a free particle, omitting all constant factors like m , 2, and \hbar . It isn't $(d^2 / dx^2)\Psi = E$, it's $(d^2 / dx^2)\Psi = E\Psi$. This fixes all the problems. For example, if we change the phase of the wavefunction by flipping its sign, the equation still holds with the same value of E . This equation is a specific example of a more general type of equation that looks like

$$\text{operator(input)} = \text{number} \times \text{input}.$$

Another, simpler example is $(d/dx)f = 3f$, which is satisfied if $f = Ae^{3x}$, where A is any constant. Such an equation says that applying the operator to the input just gives back the *input itself*, multiplied by some constant. For this reason, this type of equation is called an *eigenvalue equation*, because “eigen” is the German word for “self.” We say that 3 is the eigenvalue of the eigenvalue equation $(d/dx)f = 3f$. In the time-independent Schrödinger equation, the eigenvalue is the energy, and a solution Ψ is called a state of definite energy (or “eigenstate”).

¹“Phase” as in the phase of a wave, not as in “phase space.”

All observables in quantum mechanics are described by operators such as derivatives. The second derivative (with the appropriate factor of $-h^2/2m$) is the kinetic energy operator in quantum mechanics. Given an operator \mathcal{O} that describes a certain observable, a state Ψ with a definite value c of that observable is one for which $\mathcal{O}(\Psi) = c\Psi$. Although it's common to use parentheses when notating functions, as in $\cos(\pi) = -1$, they are optional, and we can write $\cos \pi = -1$, so we will often use notations like $\mathcal{O}\Psi$ instead of $\mathcal{O}(\Psi)$, but keep in mind that this is not multiplication, just as $\cos \pi$ doesn't mean multiplying \cos by π .

When we carried over the classical kinetic energy observable to quantum mechanics, we weren't going blind. For example, the factor of $-h^2/2m$ in front is tightly constrained by requirements like units and the need for a traveling sine wave to have positive energy. But for the superposition of two states, classical mechanics will never give us any guidance. For example, what is the body temperature of Schrödinger's cat? For the energy operators appearing in the Schrödinger equation, we used linear operators. The result was that our law of physics was perfectly linear, and this is a hard requirement, for the reasons described on p. 917. It therefore seems natural to require that *all* observables be represented by linear operators,

$$\mathcal{O}(\Psi_1 + \Psi_2) = \mathcal{O}\Psi_1 + \mathcal{O}\Psi_2.$$

Indeed, if they were not linear, then quantum mechanics would lack self-consistency, for the act of measurement can be described by applying the Schrödinger equation to a big system consisting of the system being observed interacting with the measuring device.

Finally, we have one more requirement, which is that the linear operator representing an observable should have eigenvalues that are real. This isn't because the results of a measurement must logically be real — e.g., we can measure complex impedances. But in any real-world application of the complex number system, we must always choose some arbitrary phase conventions, such as that an inductor has a positive imaginary impedance to represent the fact that the voltage leads the current by 90 degrees. (Such phase conventions are always arbitrary because we define i as $\sqrt{-1}$, but this doesn't distinguish i from $-i$.) These phase conventions are all independent of one another, and the classical ones are independent of the convention used for wavefunctions in quantum mechanics, which is that a state with positive energy twirls clockwise in the complex plane. (See also example 15, p. 992.)

Observables

In quantum mechanics, any observable is represented by a linear operator that takes a wavefunction as an input and has real eigenvalues.

Linear algebra application

Observables are represented as linear operators (p. 968). We also require that this operator have real eigenvalues.

Some important examples of observables are momentum (example 6 below), position (example 8), energy, and angular momentum. These are represented by linear operators \mathcal{O}_x , \mathcal{O}_p , \mathcal{O}_E , and \mathcal{O}_L , respectively.

The momentum operator

example 6

Quantum mechanics represents motion as a dependence of the wavefunction on position, so that a constant wavefunction has no motion. This suggests defining the momentum operator as the derivative with respect to position. This almost works, but needs to be tweaked a little. We expect that a state of definite momentum is a sine wave of the form $\Psi = e^{ikx}$. We have $k\lambda = 2\pi$ and $p = \hbar/\lambda = \hbar k$, and the sign is a matter of convention. Taking the derivative of Ψ gives an eigenvalue ik , which has the wrong units (easily fixed by tacking on a factor of \hbar), but more importantly is not real. This suggests defining the momentum operator as

$$\mathcal{O}_p = -i\hbar \frac{d}{dx}.$$

A further note about the momentum operator is example 14 on p. 988.

A nonexample

example 7

Consider the one-dimensional particle in a box, and restrict our attention to the two lowest-energy states and their superpositions. Define an operator \mathcal{O} by the rule

$$\begin{aligned}\mathcal{O}(\wedge) &= \vee \\ \mathcal{O}(\vee) &= -(\wedge).\end{aligned}$$

Since \mathcal{O} is linear, defining its action on \wedge and \vee suffices to define its action on the superpositions of these states as well. This operator has eigenvalues, one of which is i , corresponding to the state $\wedge - i\vee$. (It also has a second eigenvalue, which is imaginary as well.) Because this operator doesn't have real eigenvalues, it is not a valid observable.

Note that in examples 6 and 7, it doesn't matter whether the operator is *defined* using complex numbers. Our definition of the momentum operator was stated using an equation that had an i in it, but its eigenvalues are real, so that's OK. The operator \mathcal{O} in example 7 was defined using only real numbers, but its eigenvalues are not real.

Position is an observable

example 8

If we have a wavefunction $\Psi(x)$ expressed as a function of position x , then we simply take the operator for position \mathcal{O}_x to be multiplication by the number x ,

$$\mathcal{O}_x(\Psi) = x\Psi.$$

For example, if $\Psi = e^{ix}$ (ignoring units), then $\mathcal{O}_x(\Psi) = xe^{ix}$. This operator is definitely linear, because multiplication by a number is

linear, e.g., $7(a+b) = 7a + 7b$. The only question is whether it has eigenvalues, and whether those are real. A state of definite x , say a state with $x = 0$, would have to be represented by a wavefunction $\Psi(x)$ for which there was zero probability of having $x \neq 0$, and this requires us to have $\Psi(x) = 0$ for nonzero x . But what would be the value of $\Psi(0)$? It has to be *infinite* if Ψ is to be properly normalized. With this motivation, the physicist P.A.M. Dirac defined the Dirac delta function,

$$\delta(x) = \begin{cases} 0 & \text{for } x \neq 0 \\ +\infty & \text{for } x = 0 \end{cases}$$

Its graph is an infinitely narrow, infinitely tall spike at $x = 0$, and it has $\int_{-\infty}^{+\infty} \delta(x) dx = 1$. Mathematicians will shake their heads and say that this is not a definition of a function, but it's very useful to pretend that it is, and the delta "function" is widely used in a variety of fields such as electrical engineering. Because it was useful, mathematicians felt obliged to define a theory in which functions are generalized to things called distributions or generalized functions.

Because we represent an observable as an operator that changes a wavefunction into a new wavefunction, a common misconception is that this change represents the effect of measurement on the system. Although it is often true that microscopic systems are delicate, so that the act of measurement may have a significant effect on them, that action of the operator on the wavefunction does not represent that effect. For example, the position operator \mathcal{O}_x from example 8 consists simply of multiplication of the wavefunction by x . Suppose we have a particle in a box with a wavefunction given by $\Psi = \sin x$, where we ignore units and normalization, and the box is defined by $0 \leq x \leq \pi$. Then $\mathcal{O}_x \Psi$ eats the input wavefunction $\sin x$ and poops out the new function $x \sin x$. But the act of measuring the particle's position clearly can't do anything like this — for one thing, the function $x \sin x$ has larger values on the right side of the box than on the left, but there is nothing to create such an asymmetry in either the original state or the measuring process. The real-world effect of the measurement would probably be to knock the particle out of the box completely, since a high-resolution measurement will have a small uncertainty Δx , which by the Heisenberg uncertainty principle means creating a large Δp .

Nor is it always true that measuring a system disturbs it. For example, suppose that we prepare a beam of silver atoms, as in the Stern-Gerlach experiment, in such a way that every atom is guaranteed to be in either a state of definite $L_x = +1/2$ or $L_x = -1/2$. That is, the beam may be a mixture of both of these possibilities, but each atom is guaranteed have its spin either exactly aligned with the magnetic field or exactly antiparallel to it. Then the effect of the magnetic field is simply to sort out the two types of atoms

according to spin, without having the slightest effect on those spins.

Phase is not an observable

example 9

On p. 978 we listed three criteria for implementing the concept of an observable in quantum mechanics, and one of these was that since wavefunctions that differ only by a phase describe the same state, the result of an observation should not depend on phase. For this reason, it should not be a surprise that the mathematical definition of an observable that we came up with does not allow for the creation of an observable to describe measurement of a phase.

By way of rigorous proof, suppose to the contrary that we did have such an observable \mathcal{O}_{ph} . By our definition of an observable, it would have to have some set of eigenvalues that were real numbers. Consider such an eigenvalue φ , which might perhaps be the argument of the wavefunction in the complex plane, although we will not need to assume that. Let Ψ be the state of definite phase having the phase φ , so that

$$[1] \quad \mathcal{O}_{\text{ph}}\Psi = \varphi\Psi.$$

We can change the phase of Ψ to create a new wavefunction. Let's retard its phase by 90 degrees, creating $i\Psi$. Since Ψ was a state of definite phase, clearly $i\Psi$ is as well, and it must have some different eigenvalue φ' . Perhaps $\varphi' = \varphi + \pi/2$, but in any case we must have $\varphi' \neq \varphi$. Then

$$[2] \quad \mathcal{O}_{\text{ph}}(i\Psi) = \varphi'(i\Psi).$$

But by linearity equation [2] is equivalent to $i\mathcal{O}_{\text{ph}}\Psi = i\varphi'\Psi$, or $\mathcal{O}_{\text{ph}}\Psi = \varphi'\Psi$, and therefore by comparison with equation [1], $\varphi = \varphi'$, which is a contradiction, so we conclude that there cannot be an observable representing phase.

The result of example 9 was a bit of a foregone conclusion, since we specifically designed our notion of an observable to be insensitive to phase. Therefore this argument is subject to the objection that perhaps there is some way to measure a quantum-mechanical phase, but our definition of an observable is just too restrictive to describe it. However, we will see on p. 1001 that there are more concrete reasons why phase cannot be measured.

Time is not an observable

example 10

We do not expect to have a time operator in quantum mechanics. This follows simply because an operator is supposed to be a function that takes a wavefunction as an input, but we typically can't tell what time it is by looking at the wavefunction. For example, if the electron in a hydrogen atom is in its ground state, then we could say its energy is zero, so its frequency is zero, the period is infinite, and the wavefunction doesn't vary at all with time. (We can choose our reference level for the electrical energy U_{elec}

to be anything we like. Even if we choose it such that the energy of the ground state is nonzero, the only change in the electron's wavefunction over time will be a phase rotation, which by example 9 is not observable.)

Of course this doesn't mean that quantum mechanics forbids us from building clocks. It just tells us that many quantum mechanical systems are too simple to function as clocks. In particular, we would be misled if we pictured a hydrogen atom classically in terms of an electron traveling in a circular orbit around a proton, in which case it really could act like the hand on a tiny clock. For further discussion of this idea, see p. 1000

Since you've already studied relativity, you've had carefully inculcated in you the idea that space and time are to be treated symmetrically, as parts of a more general thing called spacetime. The differing results of examples 8 and 10 are clearly not consistent with relativity. This is to be expected because the Schrödinger equation is nonrelativistic (cf. self-check G, p. 907), and the principles laid out in this section are the principles of *nonrelativistic* quantum mechanics.

Parity

example 11

In freshman calculus you will have encountered the notion of even and odd functions. In quantum mechanics, we can have even and odd wavefunctions, and they can be distinguished from one another using the parity operator \mathcal{P} . If $\Psi(x)$ is a wavefunction, then $\mathcal{P}\Psi$ is a new wavefunction, call it Ψ' , such that $\Psi'(x) = \Psi(-x)$. In other words, the parity operator flips the wavefunction across the origin. (In three dimensions, we negate all three coordinates.) States of definite parity are represented by wavefunctions that are even (eigenvalue +1) or odd (-1).

States of definite angular momentum

example 12

In section 14.2.4, p. 964, we saw that the kinetic energy of a quantum mechanical rotor is proportional not to ℓ^2 but instead to $\ell(\ell + 1)$. This was justified qualitatively in terms of the solutions of the Schrödinger equation for a particle on a sphere, but in fact there is a deeper reason, which is that the eigenvalues of the orbital angular momentum operator turn out to be $\ell(\ell + 1)$. The parity of such a state is $(-1)^\ell$, which can be seen in figure h on p. 930.

If we have two observables, it may or may not be possible to measure them both on the same state and get exact and meaningful results. Position and momentum p and x are incompatible observables, as expressed by the Heisenberg uncertainty principle. No state is simultaneously a state of definite p and of definite x . The magnitude of an angular momentum L and its component along some axis L_z are compatible. It is common to have a state that is simultaneously a state of definite L and of definite L_z . Another example

of *incompatible* observables is L_z and L_x , as proved on p. 924.

14.6.2 The inner product

We've defined the normalization of a wavefunction as the requirement $\int_{-\infty}^{+\infty} \Psi^* \Psi dx = 1$, which means that the total probability that the particle is *somewhere* equals 1. (Another way of writing $\Psi^* \Psi$ would be $|\Psi|^2$.) This assumes that the wavefunction is written as a function of the position x . But it is also possible to have a wavefunction that depends on some other variable, such as spin or momentum, or on some combination of variables, e.g., both the spin s and the position x of an electron, $\Psi(x, s)$. We can also use a wavefunction to describe a correlation between multiple particles, in which case the wavefunction might look like $\Psi(x_1, x_2)$. The variables that the wavefunction depends on may be either continuous, like position and momentum, or discrete, like spin or angular momentum. Given all of these possibilities, we need to figure out an appropriate generalization of the integral over x that we originally used to define our normalization condition. To provide for flexibility and generality, we will start by simply defining a new notation that looks like this:

$$\langle \Psi | \Psi \rangle = 1.$$

In the case where Ψ is a function of x alone, the angle brackets $\langle \dots | \dots \rangle$ basically mean just an integral over x , and we think of the $\langle \dots |$ part as automatically implying the complex conjugation of the thing inside it. The operation $\langle \dots | \dots \rangle$ is called the *inner product*.

Because negative probabilities don't make sense, we require that the inner product of a wavefunction with itself always be positive,

$$\langle u | u \rangle \geq 0.$$

This makes it similar to the dot product used with vectors in Euclidean geometry.

In the case of Euclidean geometry, the ability to add vectors and measure their lengths automatically gives us a way to judge the similarity of two vectors. For example, if $|u| = 1$, $|v| = 1$, and $|u+v| = 2$, then we conclude that u and v are in the same direction. On the other hand, if $|u| = 1$, $|v| = 1$, and $|u+v| = \sqrt{2}$, then we can tell that u and v are perpendicular, which makes them as different as two unit-length vectors can be. More generally, $(u+v) \cdot (u+v) = |u|^2 + |v|^2 + 2u \cdot v$, because the dot product is linear, so we can see that the information about how similar u and v are is all contained in their dot product $u \cdot v$. Making the analogy with quantum mechanics, we expect that since we can define normalization of wavefunctions, we should automatically get, "for free," a way of measuring how similar two states are.

With this motivation, we assume that there is an inner product on wavefunctions that has properties analogous to those of the dot

product. We assume linearity, so that if u , v , and w are wavefunctions, then

$$\langle u|\alpha v + \beta w\rangle = \alpha\langle u|v\rangle + \beta\langle u|w\rangle$$

and

$$\langle \alpha u + \beta v|w\rangle = \alpha^*\langle u|w\rangle + \beta^*\langle v|w\rangle.$$

In the second equation, we need to take the complex conjugates α^* and β^* , for if we omitted the conjugation, then when $\langle u|u\rangle = 1$ we would have $\langle iu|iu\rangle = -1$, describing a negative probability. For similar reasons, we require that

$$\langle u|v\rangle = \langle v|u\rangle^*$$

rather than the more familiar property of the Euclidean dot product $u \cdot v = v \cdot u$.

Inner product

Wavefunctions come equipped with an inner product that has the properties described above.

If we're dealing with wavefunctions that are expressed as functions of position, then it's pretty clear how to define an appropriate inner product: $\langle u|v\rangle = \int u^*v \, dx$. The inner product axiom stated above then requires that this (possibly improper) integral converge in all cases, which means, for example, that we have to exclude infinite plane waves from consideration. However, because it's so convenient sometimes to talk about plane waves, we may break this rule when nobody is looking. Note the similarity between the expression $\int u^*v \, dx$ and the expression $u_x v_x + u_y v_y + u_z v_z$ for a dot product: the integral is a continuous sum, and the dot product is a discrete sum.

Two wavefunctions have a zero inner product if and only if they are completely distinguishable from each other by the measurement of some observable. By analogy with vectors in Euclidean space, we say that the two wavefunctions are orthogonal. For example, $\langle \curvearrowleft | \curvearrowright \rangle = 0$, as can be verified from the integral $\int_0^\pi \sin x \sin 2x \, dx = 0$. These states are also distinguishable by measuring either their momentum or their energy.

Let's consider more carefully the general justification for this assertion that perfect distinguishability is logically equivalent to a zero inner product. We have described valid observables in quantum mechanics as being represented by operators that have real eigenvalues. An alternative description of such an operator \mathcal{O} , called a hermitian

Linear algebra application

The properties listed here for inner products in quantum mechanics are just standard rules for inner products in linear algebra.

operator² after Charles Hermite, is that it is one such that for any u and v , the equation $\langle \mathcal{O}u|v\rangle = \langle u|\mathcal{O}v\rangle$ holds.³ Being hermitian is, for an operator, analogous to being real for a number. (Cf. problem 8, p. 1012.) Just as a randomly chosen complex number is unlikely to be real, a randomly chosen linear operator will almost never be hermitian. Like love, patriotism, or beauty, a nonhermitian operator fails to translate into anything a physicist can measure.

Using this alternative characterization of what makes a valid observable, we can prove, as claimed above, that if two states are distinguishable because they have definite, different values of some observable, then they are orthogonal.⁴

a / Some examples of interpretation of the inner product. The first three examples are explained immediately below. The fourth, about averages, is justified on p. 989.

$\langle \curvearrowleft \curvearrowright \rangle = 1$	The wave \curvearrowright is properly normalized.
$\langle \curvearrowleft \curvearrowleft \rangle = 0$	The waves \curvearrowleft and \curvearrowright are perfectly distinguishable.
$\langle \curvearrowright \curvearrowleft \rangle = -0.81$	The wave \curvearrowright can be expressed as $-0.81\curvearrowleft$ plus distinguishable waves. Or: measurements have probability $(-0.81)^2 \approx 0.66$ of saying one of these waves is the same as the other.
$\langle \curvearrowright \mathcal{O}_x \curvearrowleft \rangle = \frac{L}{2}$	The wave \curvearrowright has an average position $L/2$.
$\langle \curvearrowright U \curvearrowleft \rangle$	Measures the ability of an externally applied potential U to cause a jump from \curvearrowleft to \curvearrowright . The square is the transition rate per second.

Suppose that u and v are both properly normalized wavefunc-

²The mathematician's standard definition of a hermitian operator adds an additional technical condition, which is that all of the operator's eigenvalues should have magnitudes below a certain fixed bound. This is much too restrictive for our purposes, since, for example, an alpha particle in free space can have an arbitrarily large kinetic energy. In fact, nothing really bad happens if we relax our requirement for quantum-mechanical operators to be that they merely need a property called being *normal*.

³Proof that a hermitian operator has real eigenvalues: Let e be an eigenvalue, $\mathcal{O}u = eu$ for $u \neq 0$. Then $\langle \mathcal{O}u|u\rangle = \langle u|\mathcal{O}u\rangle$, so $\langle eu|u\rangle = \langle u|eu\rangle$, and $e^*\langle u|u\rangle = e\langle u|u\rangle$, so $e^* = e$, meaning that e is real.

⁴Proof: Consider states u and v with $\mathcal{O}u = e_1u$ and $\mathcal{O}v = e_2v$. If \mathcal{O} is Hermitian, we have $\langle \mathcal{O}u|v\rangle = \langle u|\mathcal{O}v\rangle$, so $e_1^*\langle u|v\rangle = e_2\langle u|v\rangle$. But since e_1 and e_2 are real and unequal, we must have $\langle u|v\rangle = 0$.

tions. If $|\langle u|v \rangle| = 1$, then the states are identical.⁵ If $\langle u|v \rangle = 0$, then u and v are completely distinguishable from one another. There is also the intermediate case where $\langle u|v \rangle$ has a magnitude greater than 0 but less than 1. In this case, we could say that u is a mixture of v plus some other state w that is distinguishable from v , i.e., that

$$|u\rangle = \alpha|v\rangle + \beta|w\rangle.$$

where $\langle v|w \rangle = 0$. We then have

$$\langle u|v \rangle = (\alpha\langle v| + \beta\langle w|)|v\rangle = \alpha.$$

Now suppose that we make measurements capable of determining whether or not the system is in the state v . If the system is prepared in state u , and we make these measurements on it, then by the linearity of the Schrödinger equation, the result is that the measuring apparatus or observer ends up in a Schrödinger's-cat state that looks like

$$\alpha|\text{observed } v\rangle + \beta|\text{observed } w\rangle.$$

We interpret squares of amplitudes as probabilities, so

$$P = |\alpha|^2 = |\langle u|v \rangle|^2$$

gives us the probability that we will have observed the state to be v . This final leap in the logic, to a probability interpretation, has felt mysterious to several generations of physicists, but recent work has clarified the situation somewhat.

On p. 940 we stated the Pauli exclusion principle by saying that two particles with half-integer spins could never occupy the same state. This was not a completely rigorous definition of the principle, since we didn't really define "same state." A more mathematically precise statement is that if one electron's wavefunction is u and another's is v , then $\langle u|v \rangle = 0$. In other words, we are ruling out not just the case where u and v are the same wavefunction, $\langle u|v \rangle = 1$, but also the intermediate case where $\langle u|v \rangle$ is greater than 0 but less than 1.

A unitary transformation is one that preserves inner products. That is, $\langle \mathcal{O}u|\mathcal{O}v \rangle = \langle u|v \rangle$. This is similar to the way in which rotations preserve dot products in Euclidean geometry. This provides a more rigorous definition of what we meant by postulating the unitary evolution of the wavefunction (p. 972). It can be shown that if the Hamiltonian is hermitian, then the evolution of the wavefunction over time is a unitary operation. This protects us from bad scenarios like the one described in example 15, p. 992.

⁵If the inner product is, for example, -1 , then the wavefunctions differ only by an unobservable difference in phase, so they really describe the same state.

Traveling waves in the quantum moat

example 13

On p. 922 we discussed the “quantum moat,” in which a particle is constrained to a circle like the moat around a castle. For the $\ell = 1$ state, the two degenerate traveling wave solutions to the Schrödinger equation are (ignoring normalization) the counterclockwise $|ccw\rangle = e^{i\theta}$ and the clockwise $|cw\rangle = e^{-i\theta}$. These states are distinguishable by their angular momenta $\ell_z = \pm 1$, so we expect them to be orthogonal. Let’s check that directly.

$$\begin{aligned}\langle ccw|cw \rangle &= \int_0^{2\pi} \left[(e^{i\theta})^* \right] e^{-i\theta} d\theta \\ &= \int_0^{2\pi} e^{-i\theta} e^{-i\theta} d\theta \\ &= \int_0^{2\pi} e^{-2i\theta} d\theta\end{aligned}$$

This is easily seen to be zero without an explicit calculation, because when we take the antiderivative of $e^{-2i\theta}$, we will get the same type of exponential, whose values when we plug in the upper and lower limits of integration will cancel each other out.

Imaginary momentum?

example 14

Here’s a paradox. If we take a wavefunction e^{rx} , where r is a constant, then applying the momentum operator $\mathcal{O}_p = -i d/dx$ (example 6, p. 980) gives

$$\mathcal{O}_p e^{rx} = -ire^{rx}.$$

For a state of definite momentum, we normally have in mind, as in examples 6 and 13, an oscillating wave where $r = ik$ is purely imaginary. But what if r is real, say $r = 1$ (ignoring units)? Then our wavefunction is e^x , and it’s a state of definite momentum — *imaginary* momentum. Oh no, what’s going on? Nice polite observables like momentum aren’t supposed to have imaginary eigenvalues.

The resolution to this paradox lies in the fundamental principles of quantum mechanics that we’ve learned. Wavefunctions are supposed to belong to a vector space in which we have a well-defined inner product. A wavefunction like $\Psi = e^x$ is ruled out by this requirement, because $\langle \Psi | \Psi \rangle$ is infinite, and therefore undefined.

Of course we could raise the same objection to a wavefunction like $\Phi = e^{ikx}$ defined for all real values of x . But when we work with wavefunctions like Φ , we usually just have in mind a computational shortcut, with the actual wavefunction being some kind of wavepacket or wave train consisting of a finite number of wavelengths. (Or we could be talking about rotation, as in the quantum moat of example 13. Note that in such an example, oscillating functions can be made to join smoothly to themselves as they wrap around, but this doesn’t work with functions like e^x .)

Averages

The average family lives down the street from me. Their family income in 2014 was \$72,641, and they have 2.5 kids. This joke depends on the fact that you can't superpose families to make a single family — but we *can* do this for wavefunctions. Suppose that the particle-in-a-box wavefunction $\hat{\psi}$ has a definite energy of 1 unit, $\mathcal{O}_E \hat{\psi} = 1\hat{\psi}$. This says that $\hat{\psi}$ is a state of definite energy 1, so that when we act on it with the energy operator \mathcal{O}_E , the result is just to multiply the wave by 1 (the eigenvalue).

If this is true, then shortening the wavelength by a factor of 2 means increasing the momentum by a factor of 2, and increasing the energy by a factor of 4. Therefore the wavefunction $\sqrt{\psi}$ has 4 units of energy $\mathcal{O}_E \sqrt{\psi} = 4\sqrt{\psi}$.

Now there is nothing wrong with mixing these together to get a state $\Psi = c\hat{\psi} + c'\sqrt{\psi}$. If both c and c' are nonzero, then we expect to get a state with properties in between those of $\hat{\psi}$ and $\sqrt{\psi}$. If we measure the energy of such a state, then our wavefunction becomes entangled with that of the particle, and we look like this:

$$c \boxed{\text{We measured the energy to be 1.}} + c' \boxed{\text{We measured the energy to be 4.}}$$

Suppose we make the mixture an equal one, $c = c'$. Then the average should be $(1 + 4)/2 = 2.5$. This turns out to be easily expressible using an inner product:

$$\langle \Psi | \mathcal{O}_E \Psi \rangle = 2.5.$$

It's a good exercise to work this out for yourself (problem 20, p. 1016). The key point is that Ψ can be expressed as a superposition of states of definite energy $\Psi = c\hat{\psi} + c'\sqrt{\psi}$, and when the operator \mathcal{O}_E works on Ψ , it gives $\mathcal{O}_E \Psi = c\hat{\psi} + 4c'\sqrt{\psi}$. (And remember that by normalization, $|c| = |c'| = 1/\sqrt{2}$.)

This is a general rule for calculating averages: for a state Ψ , the average value for an observable \mathcal{O} is $\langle \Psi | \mathcal{O} \Psi \rangle$. Because observables are hermitian, this is the same as $\langle \mathcal{O} \Psi | \Psi \rangle$.

Discussion Questions

A Suppose that by rotating vectors we could change the results of dot products. Explain why this would be very naughty, first by using an example in which $\mathbf{u} \cdot \mathbf{u} = 1$, and then, just to make it naughtier, one where $\mathbf{u} \cdot \mathbf{v} = 0$.

B Suppose that as a system evolved over time, inner products of wavefunctions could change. As in discussion question A, give shockingly naughty examples where initially we have $\langle \Psi | \Psi \rangle = 1$ and $\langle \Psi | \Phi \rangle = 0$, but later these inner products change.

14.6.3 Completeness

We have used math to back up our claim that distinguishable states are orthogonal. Going in the opposite direction, suppose that

$\langle u|v \rangle = 0$. How can we then conclude that there exists some observable \mathcal{O} that can distinguish them? There is no straightforward mathematical reason why this must be true, but it would not make sense physically to talk about two states that were utterly distinct and yet indistinguishable by any experiment. We therefore take this as a postulate.⁶

Completeness

For any system of interest, there exists a set of compatible observables, called a complete set, such that any state of the system can be expressed as a sum of wavefunctions having definite values of these observables.

The completeness postulate was discussed at a more elementary level in section 13.4.3, p. 925.

The set of wavefunctions referred to above is called a basis. (The terminology comes from linear algebra.) If we require normalization and ignore the undetectable phase, then choosing a complete set of observables is equivalent to choosing a basis. Therefore “choice of basis” and “choice of a complete set of observables” are nearly synonyms, so we will usually use the shorter phrase. Normally there is more than one possible choice of basis. The choice from among these possibilities is arbitrary, and nature doesn’t care which one we pick. That is, there is *no preferred basis*. An example of this principle is the fact that we habitually talk about “up” and “down” for the spin of an electron, which we are free to do, although it would be equally permissible to talk about left and right. Another good example is the discussion of the double degeneracy of the quantum moat on p. 922, where we were free to talk about a basis consisting of either two standing waves or two traveling waves.

As an example of the completeness principle, we have seen in the example in fig. d, p. 925, that for a rotor, the state with $\ell = 1$ and $\ell_x = 0$ can be written as a sum of the states with $\ell_z = -1$ and $\ell_z = 1$. In the language of the completeness postulate, we can express this as follows. Let our system be the set of possible states of a rotor. The observables L and L_z are compatible, and they turn

⁶Our statement of the completeness principle refers to taking a sum of wavefunctions. Because the physical motivation for the completeness postulate is so appealing, physicists are willing to stretch the definition of the word “sum” in order to make it true. The sum can be an infinite sum, and in certain cases we may even need to make it an integral, which is a kind of continuous sum. For example, consider a one-dimensional particle in a box. A complete set of observables for this system can be found by picking the energy operator alone. Now suppose we throw a particle in the box, in such a way that its position is equally likely to be anywhere in the box, i.e., its wavefunction is supposed to be constant throughout the box. Ignoring normalization, this constant wavefunction can be expressed as an infinite series in terms of the states of definite energy as $\text{---} + \frac{1}{3}\text{---} + \frac{1}{5}\text{---} + \dots$ This kind of representation of a function as an infinite sum of sine waves is called a Fourier series.

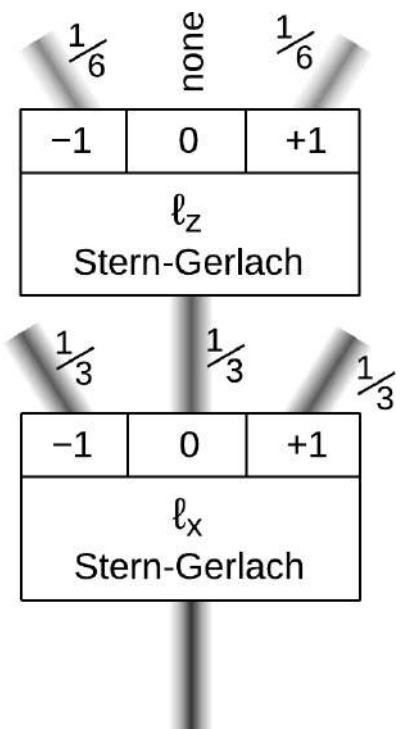
out (although we will not prove it here) to be a complete set of observables for this system. The completeness postulate is satisfied in this example because the state with $\ell_x = 0$ can be expressed as $|\ell_z = -1\rangle/\sqrt{2} + |\ell_z = 1\rangle/\sqrt{2}$.

Translating this scenario into a hypothetical real-world experiment, suppose that, as in figure b, we pass a beam of randomly oriented oxygen molecules (referred to as an unpolarized beam) through a Stern-Gerlach spectrometer that disperses them into beams with $\ell_x = -1, 0$, and $+1$. All three states are present, and in fact the beam is split into three beams of equal intensity, $1/3$ that of the original beam.⁷ Then we throw away all but the molecules having $\ell_x = 0$, and pass these through a second spectrometer, this one selecting states according to their ℓ_z . You can simulate experiments like this using an app at physics.weber.edu/schroeder/software/Spins.html. We have already found that the wavefunction of the intermediate beam is equal to the sum $|\ell_z = -1\rangle/\sqrt{2} + |\ell_z = 1\rangle/\sqrt{2}$, so interpreting squares of amplitudes as probabilities we predict a probability $(1/\sqrt{2})^2 = 1/2$ that each particle will be measured to have $\ell_z = -1$, the same probability for $+1$, and zero probability for 0. As explored in discussion question C on p. 961, this does *not* mean that the two beams that emerge from the second spectrometer have definite values of both ℓ_x and ℓ_z ; those two observables are not compatible.

In most of the examples we've encountered so far, it has been possible to think of the “wavefunction” as exactly what the word implies: a mathematical function of x (and possibly also of y , and z), whose shape we visualize as a wave. The completeness principle, however, does not assign any special role to the position operator, nor does quantum mechanics in general. And there are cases where we do not even have the option of resorting to the picture of a wave that exists in space. For example, the intrinsic angular momentum $\hbar/2$ of an electron is not a possible amount of angular momentum for a particle to generate by moving through space. In section 14.7.1, p. 993, we will discuss a very simple quantum-mechanical system consisting of an electron, at rest, surrounded by a uniform magnetic field. In this example, the motion of the electron through space is not even of interest, and a complete set of observables simply consists of L and L_z (or s and s_z , in notation that emphasizes that we're talking about intrinsic spin).

14.6.4 The Schrödinger equation in general

This raises the question of what we mean by “the Schrödinger equation” in cases where nothing is being expressed as a function of x . The basic idea of the Schrödinger equation is that a parti-



b / A beam of oxygen molecules, with $\ell = 1$, is filtered through two Stern-Gerlach spectrometers.

⁷The equality of these three intensities is not obvious geometrically, but becomes more plausible if you consider the randomness of the unpolarized beam as being defined by its having maximum entropy.

cle's energy is related to its frequency by $E = hf$, or $E = \hbar\omega$. In the form of the time-dependent Schrödinger equation that we have discussed on p. 970, $i\hbar\partial\Psi/\partial t = -(\hbar^2/2m)\nabla^2\Psi + U\Psi$, the quantity on the right-hand side of the equation is just the energy operator acting on the wavefunction. So to generalize this to cases where the wavefunction isn't expressed in terms of x , we just make that substitution:

$$i\hbar\frac{\partial\Psi}{\partial t} = \mathcal{O}_E\Psi.$$

This is as good a point as any to introduce a not-very-memorable piece of terminology, which is that the energy operator in quantum mechanics is called the *Hamiltonian*, after W.R. Hamilton. There is a classical version of the Hamiltonian, which is usually a synonym for the energy of a system, although it turns out that there are cases where it is not the same, e.g., when we adopt a rotating frame of reference. In both classical and quantum mechanics, the Hamiltonian is what determines the time-evolution of a system; in quantum mechanics, this is because it is the Hamiltonian that occurs in the Schrödinger equation. Because the Hamiltonian occurs so frequently, we will notate it as \hat{H} rather than the more cumbersome \mathcal{O}_E , where the hat is to remind us that it is an operator. A similar notation can be used for other operators when it is easier to write, e.g., \hat{s}_z rather than the clumsy \mathcal{O}_{s_z} . In the hat notation, the time-dependent Schrödinger equation looks like this:

$$i\hbar\frac{\partial\Psi}{\partial t} = \hat{H}\Psi.$$

An illegal energy operator

example 15

We have pointed out on p. 979 some reasons to think that it would be bad to have a quantum-mechanical observable whose eigenvalues were not real, i.e., one represented by a non-hermitian operator (p. 986). Even worse things happen if we try to use a non-hermitian operator for our energy operator, the Hamiltonian. As the simplest possible example, consider a system consisting of a particle at rest, and the Hamiltonian defined by

$$\hat{H}\Psi = ik\Psi,$$

where k is a nonzero real constant with units of energy. That is, the energy of the system is a constant value, which is the imaginary number ik . This operator has a single eigenvalue, ik , which is not real. The fact that it has a non-real eigenvalue is equivalent to a statement that it is non-hermitian (problem 8, p. 1012). If we plug this in to the Schrödinger equation, we get $i\hbar\partial\Psi/\partial t = ik\Psi$, or

$$\frac{\partial\Psi}{\partial t} = \frac{k}{\hbar}\Psi.$$

This differential equation is not hard to solve by the guess-and-check method. A function whose derivative is itself (except for a

multiplicative constant) is an exponential. The solution is

$$\Psi = Ae^{(k/\hbar)t},$$

where A is a constant. This is bad. Very bad. If Ψ is properly normalized at $t = 0$, then it will not be normalized at other times. If k is positive, then the total probability will become greater than 1 for $t > 0$, which we could perhaps interpret as meaning that the particle is spawning more copies of itself. Almost as bad is the case of $k < 0$, for which the particle exponentially vanishes into nothingness like the Cheshire cat. Either behavior would violate the principle of the unitary evolution of the wavefunction (p. 972).

14.6.5 Summary of the structure of quantum mechanics

We can now summarize the logical structure of quantum mechanics using the following five principles.

1. *Wavefunction fundamentalism:* All knowable information about a system is encoded in its wavefunction (ignoring phase and normalization).
2. *Inner product:* Wavefunctions come equipped with an inner product that has the properties $\langle u|\alpha v + \beta w \rangle = \alpha\langle u|v \rangle + \beta\langle u|w \rangle$ and $\langle u|v \rangle = \langle v|u \rangle^*$.
3. *Observables:* In quantum mechanics, any observable is represented by a linear operator \mathcal{O} that takes a wavefunction as an input and is hermitian, $\langle \mathcal{O}u|v \rangle = \langle u|\mathcal{O}v \rangle$.
4. *Unitary evolution of the wavefunction:* The wavefunction evolves over time, according to the Schrödinger equation $i\hbar\partial\Psi/\partial t = \hat{H}\Psi$, in a deterministic manner. Because \hat{H} is an observable, the Schrödinger equation is *linear* and also *unitary*. Unitarity means that $\langle u(t)|v(t) \rangle = \langle u(t')|v(t') \rangle$, so that probability is conserved and information is never lost.
5. *Completeness:* For any system of interest, there exists a set of compatible observables, called a complete set, such that any state of the system can be expressed as a sum of wavefunctions having definite values of these observables.

14.7 Applications to the two-state system

14.7.1 A proton in a magnetic field

As an application of the ideas discussed in section 14.6, let us consider the example of a proton at rest in a uniform magnetic field. We will find that this very simple example has surprising properties, and also that it throws light on much more general ideas than would be expected, given how specific the situation is. We discuss the

Application to MRI scans

In nuclear magnetic resonance (NMR), which is the technological basis for medical MRI scans, a very large DC magnetic field, ~ 3 T, is applied to the sample using a superconducting magnet. Protons in hydrogen atoms have their spin states split in energy by $\Delta E = 2\varepsilon = k\hbar$. After ~ 1 s, the protons reach a new thermal equilibrium state in which the probability of $|\downarrow\rangle$ and $|\uparrow\rangle$ differ by $\sim 10^{-5}$.

A brief radio-frequency pulse is then applied at the frequency ω such that $\Delta E = \hbar\omega$, so that a radio photon has the correct energy to cause a transition between the two spin states. Since there is a large number of protons, and they interact with one another, their response can be described semiclassically. The magnetization vector of the sample precesses in a complicated manner, which can be affected by the polarization and duration of the pulse.

After the radio pulse has stopped, the protons return to equilibrium again, and this changing magnetic field causes induced electric fields in a coil, which picks up a signal at the frequency ω . Spatial resolution for imaging is accomplished by adding a gradient to the magnetic field, amounting to a few percent over a distance of one meter, so that ω has different values for different points in space.

proton because the physics is then the physics of nuclear magnetic resonance (NMR), which is the technology used for, among other things, medical MRI scans.

Classically, the proton feels no magnetic force because it is at rest, and also because the field is uniform (unlike the one in the Stern-Gerlach experiment). Therefore we expect it to stay at rest. Its energy is $-\mathbf{m} \cdot \mathbf{B}$, and for the reasons discussed in sec. 11.2.4, p. 697, the magnetic dipole moment \mathbf{m} is proportional to the spin angular momentum vector \mathbf{s} , so that the energy can be broken up into a sum of three terms as $ks_x B_x + ks_y B_y + ks_z B_z$, where k is $-1/g$ times the proton's charge-to-mass ratio.

Quantum-mechanically, the components of the magnetic field will act like ordinary numbers (since the field is static, and we aren't trying to describe its dynamics quantum-mechanically), but the components of the angular momentum are observable properties of the proton, to be represented by operators. There is not always a foolproof procedure for translating a classical expression into something quantum-mechanical, but in this example it seems sensible to imagine that the classical expression for the energy can be made into a quantum-mechanical energy operator that is obtained simply by substituting the components of the angular momentum operator into the expression.

What we have determined so far is that the Hamiltonian \hat{H} will simply be a weighted sum of \hat{s}_x , \hat{s}_y , and \hat{s}_z , with the weighting determined by the components of the magnetic field.

From our previous study of angular momentum in quantum mechanics, we know that a full description of our proton's angular momentum can be given by specifying the magnitude of the angular momentum, which is a fixed $\hbar/2$, and its component along some arbitrarily chosen axis, say z . We have a state $|\uparrow\rangle$ which has eigenvalue $s_z = +\hbar/2$, and a $|\downarrow\rangle$ with $-\hbar/2$. If the magnetic field is parallel to the z axis, then the action of the Hamiltonian is easy to define in terms of these two states,

$$\begin{aligned}\hat{H}|\uparrow\rangle &= \varepsilon|\uparrow\rangle \quad \text{and} \\ \hat{H}|\downarrow\rangle &= -\varepsilon|\downarrow\rangle,\end{aligned}$$

where to keep the notation compact we write $\varepsilon = k\hbar/2$, which is an energy. The interpretation is that if there is no external magnetic field ($k = 0$), then the energies of these two states are the same (and set to zero because we choose that as an arbitrary definition), while in the presence of a B_z the two energies become unequal. The pair of states is "split" in energy by the field. Note that the above two equations are sufficient to define the Hamiltonian for *all* states, not just for states in which s_z has a definite value. This follows from the completeness principle — a state having a definite value of, say, s_x can be written as some kind of linear combination of the form

$\alpha|\uparrow\rangle + \beta|\downarrow\rangle$, and we then have $\hat{H} = \alpha\varepsilon|\uparrow\rangle - \beta\varepsilon|\downarrow\rangle$.

Now suppose that the magnetic field is not parallel to the z axis. One way to handle this situation would be simply to redefine the coordinate system so that the z axis was back in alignment with the direction of the field. But suppose that's not convenient. Then the Hamiltonian will have a different form. But because the Hamiltonian must be Hermitian (see p. 986), there is not much freedom in choosing this form. It must look something like this:

$$\begin{aligned}\hat{H}|\uparrow\rangle &= \varepsilon|\uparrow\rangle + f|\downarrow\rangle \\ \hat{H}|\downarrow\rangle &= f^*|\uparrow\rangle - \varepsilon|\downarrow\rangle.\end{aligned}$$

Here the constant f is a complex number with units of energy. The interpretation is that ε tells us how much energy splitting we would have had if the magnetic field had not had any x or y components, while f brings in the effect of those components. We could go ahead and work out the eigenvalues of this operator by writing down the eigenvalue equation and solving it by brute force, but the result is likely to seem less mysterious if we instead apply the following physical argument.

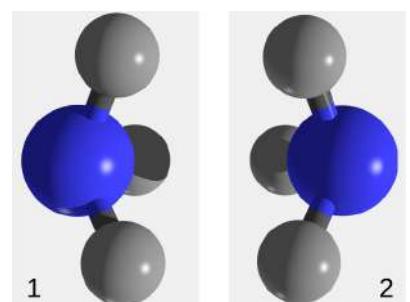
Although f lies at some point in the complex plane with some phase angle $\arg f$, such phase angles in quantum mechanics are not directly observable. Since energies *are* observable, it follows that the two eigenvalues of energy can only depend on the magnitude of f , not on its phase. By rotational invariance (sec. 3.4.2, p. 195), we also know that these energies can only depend on $|\mathbf{B}| = \sqrt{B_x^2 + B_y^2 + B_z^2}$, and in fact when the direction of the field is fixed they must be proportional to $|\mathbf{B}|$ (not to, e.g., the cube of the field). We have already interpreted ε as being essentially B_z , except for a constant of proportionality, so it follows from units that the energies must be of the form $E = \pm\sqrt{\varepsilon^2 + (\dots)|f|^2} = \pm\sqrt{\varepsilon^2 + (\dots)f^*f}$, where (\dots) represents a universal unitless constant, which turns out to be 1. We therefore have for the energies the result

$$E = \pm\sqrt{\varepsilon^2 + f^*f}.$$

Note that our earlier result of $E = \pm\varepsilon$ is recovered when $f = 0$.

14.7.2 The ammonia molecule

I chose the example of the proton in a magnetic field in the preceding section for ease of computation, but the treatment of the general case where $f \neq 0$ may not have seemed especially compelling, since we would always have the freedom to align our z axis with the field, giving $f = 0$. But our results from that analysis are of much greater generality. They do not depend on any facts about the system other than the fact that it is a system with two states. To see the full power and generality of this approach, we will apply it to the ammonia molecule, NH_3 , shown in figure a.



a / The ammonia molecule, in states that are inverted relative to one another.

At ordinary temperatures, this molecule is likely to be rotating, and its angular momentum will have some component about its symmetry axis (the left-right axis in the diagram). Let's say, for example, that the angular momentum vector points to the right, which we'll say is the positive x direction. Then the two orientations of the molecule shown in figure a are distinguishable. In one, the electric dipole vector (example 6, p. 588) points in the same direction as the angular momentum vector, and in the other they point in opposite directions.⁸ For a fixed angular momentum, we have a two-state system, as in section 14.7.1.

Classically, the molecule's moment of inertia is the same for orientations $a/1$ and $a/2$, so we would expect there to be two states with the same energy. We can always add an arbitrary constant to the energies, so if they're the same, we can just say they're both zero. Does this mean that quantum-mechanically, we simply have $\hat{H} = 0$? That would be boring. But this cannot be true, for the following reason. According to the Schrödinger equation, a state of definite energy is a state that has a definite frequency, so it lasts forever, just twirling its phase angle around in the complex plane at a rate $\omega = E/\hbar$. So if state 1 were a state of definite energy, then according to the Schrödinger equation if we initially put the molecule in state 1 it would stay in that state forever. But this cannot be the case, because we know it is possible for the molecule to switch from state 1 to state 2 by turning itself inside out like an umbrella caught by a gust of wind. The possibility of this type of inversion is not just an optional thing. Vibrations that flex the shape will exist due to zero-point motion (p. 965). Even if inversion requires a lot of energy, and the molecule doesn't have that much energy, there is at least some probability of having quantum-mechanical tunneling from 1 to 2. If we prepare the molecule in state 1, and then observe it at some later time, there is some nonzero probability of finding it in state 2. This is a contradiction, so our assumption of $\hat{H} = 0$ must have been false.

So the Hamiltonian is not zero, but we already know the full variety of forms that the Hamiltonian of a two-state system can have. We only have a couple of parameters to play with, the numbers ε and f . We have $\varepsilon = 0$ by symmetry, so the only possible form for the Hamiltonian is this:

$$\begin{aligned}\hat{H}|1\rangle &= f|2\rangle \\ \hat{H}|2\rangle &= f^*|1\rangle.\end{aligned}$$

⁸This argument shows that when $L_z \neq 0$ we have two distinguishable states, but it does not necessarily tell us anything about the converse. When $L_z = 0$, are there two states, or only one? The analysis in this case is rather intricate, and depends on the Pauli exclusion principle and the fact that the hydrogen atoms are all identical, that there are three of them, and that their nuclei are fermions. See Townes and Schawlow, *Microwave Spectroscopy*, 1955, pp. 69-71.

Because we can define the states $|1\rangle$ and $|2\rangle$ with any phases we like, we are free to take f to be real, $f^* = f$, although this implies a certain relationship *between* the phases of $|1\rangle$ and $|2\rangle$. If we visualize these states as bell-shaped functions of an x coordinate describing the position of the nitrogen relative to the plane of the hydrogens, then it would be nice to have a phase convention such that where the tails of the wavefunctions overlap, inside the barrier, they have the same phase. This turns out to be the case when f is real and negative, so we will assume that from now on. Recycling our previous result for the energies, we have $E = \pm\sqrt{\varepsilon^2 + f^2} = \pm f$. If the tunneling probability approaches zero, then we expect f to go to zero, and the energy splitting approaches zero, as we had expected classically. Experimentally, we do observe these two states in ammonia. The difference in energy is extremely small — e.g., for the state with angular momentum $1\hbar$ it is about 9.8×10^{-5} eV, so that if a photon is emitted or absorbed in a transition between the states, it lies in the microwave spectrum. This energy difference equals $2|f|$, and its smallness indicates that the tunneling probability is small.

Let's find the states of definite energy for this system. For the ground state, whose energy is $-|f|$, we need to look for a state of the form $|g.s.\rangle = (\dots)|1\rangle + (\dots)|2\rangle$ such that $\hat{H}|g.s.\rangle = -|f||g.s.\rangle = f|g.s.\rangle$. If we don't worry about normalization or an over-all phase, we are free to take the first (...) equal to 1, so that $|g.s.\rangle = |1\rangle + \alpha|2\rangle$, for some complex number α . We then have

$$\begin{aligned}\hat{H}|g.s.\rangle &= \hat{H}(|1\rangle + \alpha|2\rangle) \\ &= f|2\rangle + \alpha f|1\rangle,\end{aligned}$$

and setting this equal to $f|g.s.\rangle$ gives $\alpha = 1$, so that

$$|g.s.\rangle = |1\rangle + |2\rangle.$$

The coefficients (...) that we set out to find are both equal to +1. Their equal magnitudes tell us that the ground state is one in which the molecule has an *equal* probability of existing in either inversion. Since the two coefficients are both positive, and we have defined $|1\rangle$ and $|2\rangle$ such that their phases agree when they overlap inside the barrier, this is a state of positive parity. The determination of the excited state is left as an exercise, problem 10 on p. 1013.

From a classical point of view, we would think of the set of states

$$\{ |1\rangle, |2\rangle \}$$

as the natural way of describing the possible states of the system. These two states are the ones that we can draw pictures of, a/1 and a/2. But part of the structure of quantum mechanics is that there is *no preferred basis* (p. 990), and there is nothing wrong with using the ground state and first excited state to form the basis

$$\{ |g.s.\rangle, |ex.s.\rangle \}$$

instead. In the language of the completeness principle (p. 990), one possible choice of a complete set of compatible observables for this molecule is the set consisting of a single observable, the energy. The {ground-state, excited-state} basis just happens to be the one associated with this particular observable. If the ammonia molecule had just broken off from some larger molecule, then it would be oriented in a specific direction, and we would probably find it more convenient to describe it in the {1,2} basis.

14.8 Energy-time uncertainty

14.8.1 Classical uncertainty relations

Consider the following classical system of analogies.

$$\begin{array}{lll} \text{space} & x & k \\ \text{time} & t & \omega \end{array} \quad \Delta x \Delta k \gtrsim 1 \quad \Delta t \Delta \omega \gtrsim 1$$

Here the quantity $k = 2\pi/\lambda$ is called the wavenumber. The inequality $\Delta x \Delta k \gtrsim 1$ is a kind of classical uncertainty relation that is closely related to the Heisenberg uncertainty principle. Its classical nature is immediately apparent because it doesn't involve Planck's constant. If you look back at the argument given on p. 903 to justify the Heisenberg uncertainty principle, you will see that it carries through equally well if we simply omit the quantum-mechanical ingredients and use it to put a bound on $\Delta x \Delta k$ instead of $\Delta x \Delta p$. Once we've established the bound on $\Delta x \Delta k$, the one on $\Delta x \Delta p$ follows immediately because $p = h/\lambda = \hbar k$.

The second line of the table is in strict analogy to the first line. A good practical example is the high-speed transmission of digital data over transmission lines such as fiber-optic cables. Suppose that we wish to send a string of 0's and 1's, and a 1 is to be represented by a square pulse. If we want to transmit the data at high speed, then we need the duration Δt of this pulse to be short, perhaps in the microsecond or even nanosecond range. This cannot be done if the signal consists only of a single frequency. A signal that only contains a single, pure frequency is just a sinusoidal wave that has existed infinitely far back in the past and will exist infinitely far into the future. Such a wave carries no information at all. Our frequency-time uncertainty relation tells us that if the duration of a pulse is to be, say, a microsecond, then the signal's spread in frequency must be at least on the order of 1 MHz. This is why we use the term "bandwidth" to describe the speed of a communication channel.

14.8.2 Energy-time uncertainty

In a quantum-mechanical context, we have $E = \hbar\omega$, so there is an energy-time uncertainty relation,

$$\Delta E \Delta t \gtrsim \hbar.$$

As with the Heisenberg uncertainty principle for momentum and position, the symbol \gtrsim means that we leave out a numerical factor, which can only be precisely defined if we fix some specific statistical definition of Δ , e.g., a standard deviation.

The interpretation of the energy-time uncertainty relation is a little tricky, because although the classical analogy between space and time is exact, the quantum-mechanical analogy breaks down. This is because time in nonrelativistic quantum mechanics, unlike position, is not an observable (example 10). Time in this theory is just a universal parameter. The physicist Lev Landau liked to tell his students that there was no energy-time uncertainty relation, because “I can measure the energy, and look at my watch; then I know both energy and time!” One good way of interpreting it is that if there is a transfer of energy between two systems, then it relates the uncertainty ΔE in the amount of energy transferred during the duration Δt of the interaction.

For example, suppose we wish to bounce a photon off of a hydrogen atom in order to determine whether the atom is in its ground state. This is not necessarily an easy thing to do by extracting whatever information we get from the reflected photon, but the ground state is orthogonal to the other states, so we are at least encouraged to believe that it is not theoretically impossible. But there is a hard theoretical limit on how *quickly* we can make such a determination. The difference in energy between the ground state and the first excited state is 1.6×10^{-18} J, so we must use a photon with an energy less than this amount, or else the act of observing the atom may in fact destroy the property we were hoping to measure. By the energy-time uncertainty relation, this implies that the measurement process cannot be done in less than about 10^{-15} seconds. This example may seem impractical, but in fact computer memories are starting to reach the level of speed and miniaturization at which such fundamental constraints become relevant.

Mortality for hydrogen

example 16

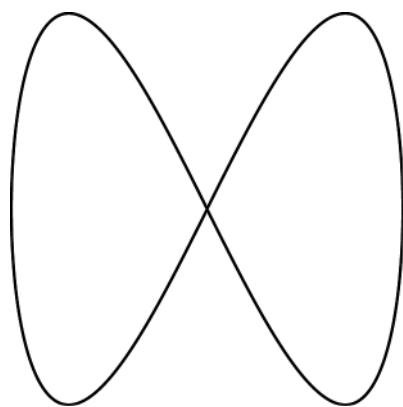
In atomic physics, when a photon is emitted or absorbed it is almost always in a wave pattern with angular momentum 1 (i.e., $1\hbar$) and negative parity (example 11). Classically, this is the type of radiation pattern that we would get from an electric dipole spinning end over end, so we call it an electric dipole transition. Because the electromagnetic interaction has a symmetry between left- and right-handedness (section 11.1.5, p. 687), this means that an electric dipole transition can never cause a transition from one state of an atom to another state with the same parity.

Now the ground state of the hydrogen atom has $\ell = 0$ and is therefore a state of positive parity. One of the first excited states, referred to as the 2s state, also has these properties, and therefore it is impossible for the 2s state to decay to the ground state

by emitting an electric dipole photon. The happy atom probably believes that once it's in the exalted 2s state, it can stay that way forever. One way for it to be cheated of immortality is if it undergoes a collision with another atom, but in some so-called planetary nebulae (hot clouds of gas cast off by dying stars), the density can be so low that collisions are very infrequent. In this situation, the dominant process for decay of the 2s state can be the simultaneous emission of *two* photons. An exact and rigorous calculation of the rate of decay for this process is quite technical, but a fairly reasonable estimate can be obtained by the following semiclassical argument based on the energy-time uncertainty relation.

The typical rate of emission for a photon, when not forbidden by parity, is $R \sim 10^9 \text{ s}^{-1}$, i.e., it takes about a nanosecond. We can think of the two-photon decay as an energy-nonconserving jump up to some *higher-energy* state, with the emission of a photon, followed by the emission of a second photon leading down to the ground state. The first jump can happen because of the energy-time uncertainty relation, which allows the electron to stay in the intermediate state for a time $t \sim h/E$, which is on the order of 10^{-15} s . The probability for the second photon to be emitted within this time is Rt , so the rate for the whole two-photon process is $R^2 t \sim 10 \text{ s}^{-1}$. Considering the extremely crude nature of this calculation, the result is in good agreement with the observed rate of about 0.1 s^{-1} . The process is actually observed, and contributes a continuous background spectrum in addition to the discrete line spectrum when such nebulae are observed with a spectrometer through a telescope.

A fundamental application of the energy-time uncertainty relation is to the explication of what it means to measure time in quantum mechanics. In example 10 on p. 982 we argued that time is not an observable in quantum mechanics because time cannot in general be measured by looking at a quantum-mechanical system: many quantum-mechanical systems are too simple to function as clocks. We can now see in more detail what “too simple” might mean here. Microscopic systems, unlike macroscopic ones, are often encountered in a definite state of energy, such as the ground state. Such a state has $\Delta E = 0$ and therefore by the energy-time uncertainty relation it has $\Delta t = \infty$. In other words, the only time evolution in such a system consists of the system's over-all phase twirling in the complex plane at a steady rate, but phase isn't measurable, so we can't use this rotation like the hand on a clock. To make a clock, we need, at a bare minimum, a system that is in a superposition of *two* different energy levels. We then have two independent phases. Although absolute phases are not measurable, relative ones are, and for example when we measure a double-slit interference pattern, that is exactly what we are doing: observing



a / The Lissajous figure $x = \cos t$,
 $y = \sin 2t$.

(statistically) the difference between two phases. As a loose conceptual analogy, this is like the idea that a figure-eight Lissajous pattern has an identifiable feature where it crosses itself, the crossing being like the tick of a clock.

14.9 Randomization of phase

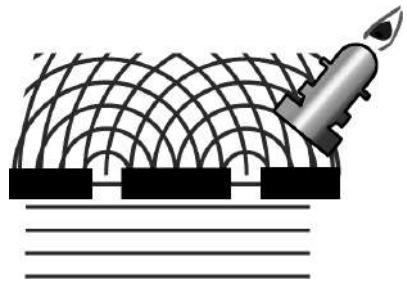
14.9.1 Randomization of phase in a measurement

The energy-time uncertainty relation can help us to understand one of the most puzzling issues in quantum mechanics, which is the problem of measurement. What happens when we use a macroscopic measuring device, which is well described by classical physics, to observe a microscopic system, which is quantum-mechanical? How do we reconcile these two seemingly incompatible descriptions of reality when both appear to be in play simultaneously?

Consider an electron passing through a double-slit apparatus. We have already considered the possibility of covering one slit (discussion question D, p. 883). Suppose instead that we carefully watch one slit through a microscope, and see whether or not the electron passed through it. If we could perform this observation without disturbing the electron, then a paradox would arise. For if we haven't disturbed the electron, then there should still be a double-slit interference pattern. But if we watch one slit, then half of the time we should see that the electron did not go through it, and therefore the slit's existence is of no importance, and we can't possibly get a double-slit interference pattern.

To avoid this contradiction, it appears that nature must conspire against us in such a way that observing the slit inevitably *does* disturb the electron. The energy-time uncertainty relation explains why this is so. Our observation of the electron is an interaction between the electron and our macroscopic measuring device. This interaction will presumably transfer some amount of energy E into or out of the electron, and if our goal was to avoid disturbing the electron, we would imagine that it would be best to make E very small. But the energy-time uncertainty $\Delta E \Delta t \gtrsim h$ relation tells us that if this energy is to have a value that is confined to some small range ΔE , then the time Δt it takes for the interaction to occur must be at least $\sim h/\Delta E$. While the electron is being subjected to this interaction, its phase is rotating around the complex plane like $e^{i\omega t} = e^{iEt/\hbar}$. The total change in the phase angle $\phi = E\Delta t/\hbar$ is uncertain because E is uncertain, so our observation will inevitably change the phase by some random amount, which is uncertain by an amount $\Delta\phi = \Delta E \Delta t / \hbar$, so $\Delta\phi \gtrsim 1$.

Thus is won't actually help us if we make the interaction very gentle, because the lengthening of the time has a compensating effect. Any slight alteration in the frequency will have more time to



a / Spying on one slit in the double-slit experiment.

accumulate into a big phase difference, and we still end up with a phase uncertainty that is at least on the order of 1. Although we haven't stated our uncertainty relations with enough mathematical precision to state this lower bound with all the right factors of 2 and π , it turns out that $\Delta\phi \geq 2\pi$. That is, any such observation will have the effect of *completely* randomizing the phase of the thing being observed. In fact, macroscopic measuring devices normally exceed the bounds set by the uncertainty relations by many orders of magnitude, so there will typically be a vast amount of overkill in this randomization. This is a general rule for reasoning about quantum-mechanical measurements: they always completely randomize the quantum-mechanical phase of the thing being measured. This provides a more physical justification for our more abstract mathematical proof in example 9 on p. 982 that phase is not an observable.

In our example of the double slit, what will be the effect of this randomization of the electron's phase? In our usual description of the double slit, we assume that the circular waves emerging from one slit are in phase with those that come out through the other one, so that the double-slit interference pattern has a maximum in the center. But if, for example, one of the waves has its phase inverted, then all the maxima of our interference pattern will become minima and vice versa. If the phase is randomized, then the positions of the maxima and minima are randomized as well, and thus if we try to collect data on enough electrons to see an interference pattern, we will not see maxima and minima at all.

One subtle question about this description is the following. The randomization of the phase by the measurement appears to have erased the information about the phase relationship between the parts of the wave in the two slits. But how can this be, since one of our principles of quantum mechanics (p. 993) is that time evolution is always unitary, so no information is ever supposed to be lost? The resolution of this paradox is that the phase information still exists, but it has been taken away from the electron and flowed out into the observer and the environment. This is similar to the classical paradox of what happens to the (classical) information written on a piece of paper when we burn the paper: the information still exists, and could in principle be reconstructed by observing all the molecules and tracing their trajectories back in time using Newton's laws.

14.9.2 Decoherence

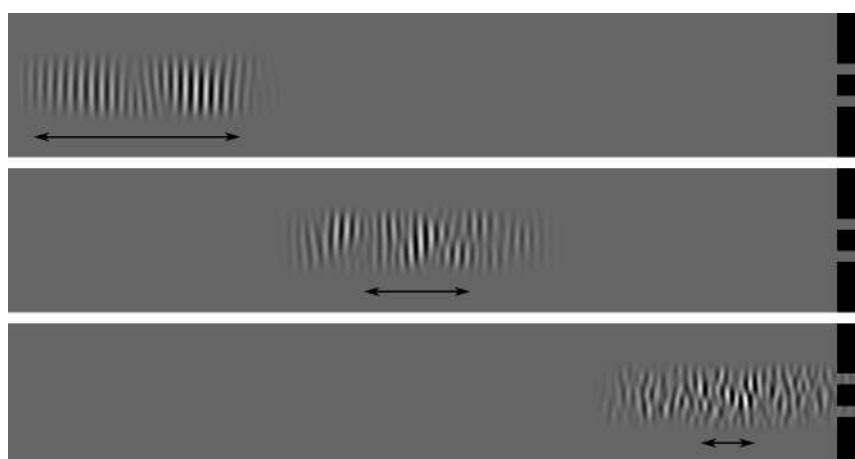
Starting around 1970, physicists began to realize that ideas involving a loss of coherence, or "decoherence," could help to explain some things about quantum mechanics that had previously seemed mysterious. The classical notions of coherence and coherence length were described in sec. 12.5.8, p. 825, and quantum-mechanical de-

coherence was briefly introduced on p. 887.

One mystery was the fact that it is difficult to demonstrate wave interference effects with large objects. This is partly because the wavelength $\lambda = h/p = h/mv$ tends to be small for an object with a large mass. But even taking this into account, we do not seem to have much luck observing, for example, double-slit diffraction of very large molecules, even when we use slits with appropriate dimensions and a detector with a good enough angular resolution.

In the early days of quantum mechanics, people like Bohr and Heisenberg imagined that there was simply a clear division between the macroscopic and microscopic worlds. Big things and small things just had different rules: Newton's laws in one case, quantum mechanics in the other. But this is no longer a tenable position, because we now know that there is no limit on the distance scales over which quantum-mechanical behavior can occur. For example, a communication satellite carried out a demonstration in 2017 in which a coherence length of 1200 km was demonstrated using photons.⁹

The insight about decoherence was the following. Consider the most massive material object that has so far been successfully diffracted through a grating, which was a molecule consisting of about 810 atoms in an experiment by Eibenberger *et al.* in 2013.¹⁰ While this molecule was propagating through the apparatus as a wave, the experimenters needed to keep it from simply being stopped by a collision with an air molecule. For this reason, they had to do the experiment inside a vacuum chamber, with an extremely good vacuum. But even then, the molecule was being bombarded by photons of infrared light emitted from the walls of the chamber. The effect of this bombardment is to disrupt the molecule's wavefunction and reduce its coherence length (p. 826).



b / A large molecule such as the one in the Eibenberger experiment is represented by its wavepacket. As the molecule starts out, its coherence length, shown by the arrows, is quite long. As it flies to the right, it is bombarded by infrared photons, which randomize its phase, causing its coherence length to shorten exponentially: by a factor of two in the second panel, and by a further factor of two in the final one. When the packet enters the double slit, its coherence length is on the same order of magnitude as the slits' spacing d , which will worsen but not entirely eliminate the observability of interference fringes. (This is only a schematic representation, with the wavepacket shown as being many orders of magnitude bigger than its actual size in relation to

⁹Yin *et al.*, arxiv.org/abs/1707.01339

¹⁰arxiv.org/abs/1310.8343

This causes an effect similar to the one in the situation illustrated in figure a, where we spy on one slit of a double-slit apparatus. The microscope would operate by bouncing photons off of the electron, and the result is to disrupt the coherence of the electron's wavefunction, so that the coherence length is no longer as large as the distance between the slits. The infrared photons in the Eibenberger experiment were not introduced intentionally, but they were still bouncing off of the molecules and producing a similar decrease in the coherence length. This decoherence effect was the reason that the experiment was limited to molecules of the size they used. Even though the molecules took only about 400 nanoseconds to fly through the apparatus, there was a significant amount of decoherence. A larger molecule would have been a bigger target for photons and would have undergone decoherence more quickly, making interference unobservable.

As in the example of spying on one slit of a double-slit experiment, the question arises of what has happened to the phase information that appears to have been erased by decoherence, violating unitarity. The resolution is the same (p. 1002): the information has flowed out into the environment, but is no longer in a form in which it is practical to recover it.

14.10 Quantum computing and the no-cloning theorem

Computers and information transmission systems such as the internet are currently implemented as classical devices. For example, the wavelengths of the electrons that carry signals in a computer chip are currently orders of magnitude shorter than the size of the logic gates, so that wave effects such as diffraction and interference are not important (problem 22, p. 946). Even if the current devices such as silicon chips and fiber-optic cables could simply be scaled down to sizes comparable to the electrons' wavelengths, quantum effects would at some point simply make them start breaking down or behaving unreliably.

It is possible, however, to design qualitatively different devices in which information and signals are intentionally manipulated in an explicitly quantum-mechanical fashion. This is the frontier known as quantum computing. In a quantum computer, the basic unit of information is not the classical bit but the quantum bit or *qubit*. A qubit can exist in a superposition of the 0 and 1 states, with a well defined phase, e.g., $\Psi_0 + \Psi_1$ is a different state than $\Psi_0 - \Psi_1$ or $\Psi_0 + i\Psi_1$. Furthermore, one qubit can have its state entangled with another's. For example, $\Psi_{01} + \Psi_{10}$ describes a state in which we have two bits, neither in a definite 0 or 1 state, but which are guaranteed to add up to 1. That is, if one is true, then the other is guaranteed to be false. It has been shown that some problems that are hard

for classical computers are more tractable for a quantum computer. For example, there is a known quantum computing algorithm that is capable of efficiently factoring large integers, and when this is eventually implemented in a practical device, it will have the effect of breaking the cryptographic algorithms that you currently use for online privacy and security, since the security of those algorithms is predicated on the assumption that factorization is hard. This would be a disaster for online economic activity and could have effects such as unmasking political dissidents.

A different application, and one that is easier to explain, is that quantum computing makes it possible in theory to make copy-proof information. This would not be useful to Hollywood studios trying to prevent copying of their movies, since the images have to pass through classical devices anyway in order to be displayed, but it means that one might be able to send private information through a quantum internet in such a way that it could not be copied by snoops, even in theory. In contrast, current classical methods of encryption are designed to allow eavesdropping on an information packet as it hops across the internet, but to make the copy useless to prying eyes because it cannot be decoded.

The theoretical key to this application of quantum computing is the counterintuitive *no-cloning theorem*, which states that it is not possible to make a copy of an unknown quantum state.¹¹ To see why this works, suppose that we implement a qubit using the spin 1/2 of a silver atom, with the convention that the 0 state is represented by $s_x = -1/2$ and 1 by $s_x = +1/2$. If you provide me with an atom that you have prepared, then it might seem straightforward, at least in principle, for me to copy its state. I can shoot it through a magnetic spectrometer, as in the Stern-Gerlach experiment, and measure its s_x . Then I prepare another silver atom in the same state. What's the problem?

The problem is that if the state of the atom is truly unknown to me, then I have no way of knowing that it is actually in one of the two states $s_x = -1/2$ and $s_x = +1/2$. It could instead be in some superposition of these, such as $\Psi_{s_x=-1/2}/\sqrt{2} + i\Psi_{s_x=+1/2}/\sqrt{2}$, with a 90-degree phase angle between the two components. Then when I send your atom through the spectrometer, the world becomes one in which both the spectrometer and my brain are in a superposition of the two states. In one of these worlds, I then go ahead and prepare my copy-atom in the $s_x = -1/2$ state, and in the other one I set it up as a $+1/2$. You could say that my copy-atom is, like the original, in a superposition of the two s_x states, but there is no reason to think that it will be the *same* superposition, with the same 90-degree phase angle. In fact, by the argument on p. 1001 we know

¹¹The prohibition actually only applies to making a copy that can be separated from the original. For the complete statement of this, see p. 1009.

that it is not possible by *any* measurement to extract this phase information and convert it into classical information. Furthermore, our final result is not really as simple as a copy-atom in some unknown superposition of the two s_x states. It is a silver atom whose spin is correlated with the state of the original, but also correlated with the state of the spectrometer and the state of my brain.

The impossibility of copying an unknown quantum state is enforced by nature in full generality, not just by the specific mechanisms described in the artificial scenario described above. To see why, consider what would happen to the state of the “blank” atom on which we had hoped to impose the copied state. Its state would have been overwritten, but this would imply a loss of information, which is forbidden by the unitarity postulate of quantum mechanics (p. 972).

The no-cloning theorem would seem to severely limit the practicability of quantum computing. When you run a program on a classical computer, the very first step to be performed by the operating system is to copy the program’s code and data from storage into random-access memory. If a quantum computer can’t copy anything, then how do we perform this initial step? But the no-cloning theorem doesn’t actually forbid copying *any* quantum state — it forbids copying an *unknown* state. Going back to the example of the silver atom, imagine that rather than presenting me with a silver atom in a completely unknown quantum state, you give me a solemn promise that it will be either in the state $s_x = -1/2$ or the state $s_x = +1/2$ — not some superposition of these. Then if you trace back through the logic of the scenario, you will find that there is absolutely nothing preventing me from making an accurate copy.

Once the software on a quantum computer starts running, its qubits will certainly start going into superpositions of the 0 state and the 1 state. By the no-cloning theorem, these cannot be copied from one memory location to another, overwriting the previous contents of the target location. But that simply isn’t how quantum computing works. Rather than attempting to copy, erase, and overwrite bits as in a classical computer, the software is designed to create complicated correlations between the different bits. This model of computing is not necessarily better or worse over all than classical digital computing, but it differs from it as much as an iPhone’s model of computing differs from that of a slide rule.

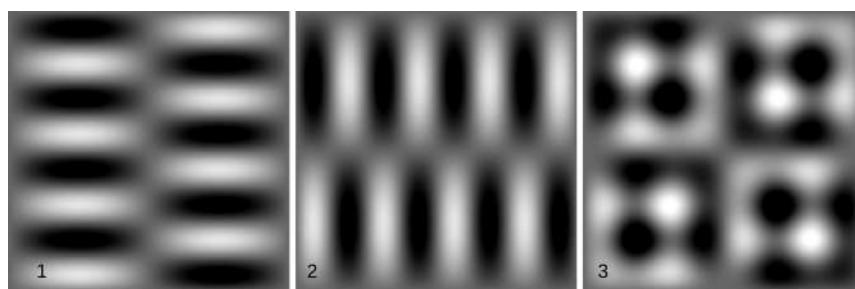
When a classical computer such as a cash register or phone is done with its computation, we have to find out the result through an output such as a paper tape or LCD screen. These are classical devices. If a quantum computer is to produce a result for use by humans, then it will also need to send its output through a classical device. We might hope to be able to convert the quantum information faithfully into classical information. But we can prove

based on the no-cloning theorem that such a conversion will always be “lossy” — will always involve a degradation of the information. A lossless conversion, such as a unit conversion, is one that can be done as a round-trip, e.g., $1\text{ m} \rightarrow 100\text{ cm} \rightarrow 1\text{ m}$, with the final result being identical to the original. If we could completely encode qbits into bits, then we could make a second copy of the bits and violate no-cloning by converting back to qbits. This is a contradiction, so we conclude that lossless conversion of classical information to quantum information is impossible.

14.11 More about entanglement

A basic difference between classical computing and quantum computing is that qubits can be entangled with each other. We’ve only discussed entanglement briefly in sec. 13.2.4, p. 884, where the basic idea was that either Alice or Bob could detect a certain photon, but not both. Alice and Bob’s states were entangled, as were the macroscopic diamonds in the 2012 real-world experiment described on p. 888. More generally, what is entanglement?

Entanglement is the opposite of separability (sec. 14.5.2). To see what is meant by this statement, consider figure a. In a/1, we have the function $\Psi_1 = \sin x \sin 4y$. This could be a two-dimensional particle in a box, with a certain amount of momentum in the x direction, and four times that momentum in the y direction. It is because the Schrödinger equation for the particle in a box is separable in x and y that we can write down this wavefunction by multiplying two different one-dimensional wavefunctions. In figure a/2, Ψ_2 is like Ψ_1 but with x and y interchanged, while a/3 shows the superposition $\Psi_3 = (\Psi_1 + \Psi_2)/\sqrt{2}$.



a / States of a particle in a box that are separable in terms of p_x and p_y (1 and 2) and entangled (3, a superposition of 1 and 2).

From a fancier theoretical point of view, we could say that this system, which seems like a single thing (the particle), is actually built out of two subsystems. One subsystem is the motion in the x direction, and the other is the y . The fact that the Schrödinger equation is separable can be interpreted as being because the x and y motion are independent of one another. Exactly the same thing would happen if this were a classical pool ball on a square table. Its

x and y motion don't affect each other, and, e.g., if the ball hits the right-hand cushion and has its x momentum reflected, that doesn't change its y momentum. It's as if the pool ball in two dimensions were really two different beads, one sliding along a wire parallel to the x axis and the other sliding up and down. In either the classical case or the quantum-mechanical case, we have built a composite system out of two independent subsystems.

In an example like Ψ_1 , it is possible to assign a definite state to the subsystems: continuing to ignore units, we can write $p_x = \pm 1$ and $p_y = \pm 4$. The state with wavefunction Ψ_2 has the same energy as Ψ_1 , and again the subsystems have a definite state, $p_x = \pm 4$ and $p_y = \pm 1$.

But for the superposition Ψ_3 , this is no longer true. If we measure either p_x or p_y for this state, we may get either ± 1 or ± 4 , with equal probability. We say that this state is entangled in the same way that Alice and Bob were entangled on p. 885. Neither Bob nor Alice is in a definite state of I-saw-a-photon or I-never-saw-a-photon. However, if we ask Bob whether he saw a photon, and he says yes, then we gain information about Alice: that she didn't see a photon. Similarly, if we measure p_x for the particle in state Ψ_3 and get -4 , then we gain information about p_y : we know that it is ± 1 .

Because separable states are the simplest things we can make by putting together subsystems like legos, it's convenient to have a notation for them. In the angle-bracket notation, all of the following are possible ways that people might notate a state like Ψ_1 :

$$|1, 4\rangle \quad \text{or} \quad |1\rangle |4\rangle \quad \text{or} \quad |1\rangle \otimes |4\rangle$$

The cross with a circle around it, \otimes , doesn't really indicate multiplication. It's more like a punctuation mark or a conjunction, meaning "and also," as in, "I'll have the eggplant, and also a beer." It's called a tensor product, which makes it sound scary.

To show the generality of the idea of entanglement, let's consider an example from particle physics. The π^0 is a particle that participates in strong nuclear interactions, and therefore can be created in nuclear reactions. It's known as a pion. There are other types of pions. The π^0 is the only electrically neutral one, hence the superscript 0. All pions are unstable, which is why we need to create them in reactions rather than looking for them in rocks and trees. The π^0 has a half-life of only 10^{-16} s, and one of the ways in which it can decay is into an electron and a positron (antielectron),

$$\pi^0 \rightarrow e^- + e^+.$$

You can verify that charge is conserved in this reaction. In the frame of reference where the pion is initially at rest, the speeds of the electron and positron are fixed by conservation of energy and momentum, so there is not much that is interesting to measure about

them other than their spins. The pion has zero spin, which makes it somewhat unusual in the world of particle physics. If we assume as well, for simplicity, that the electron and positron don't have any orbital angular momentum, then by conservation of angular momentum, the spin-1/2 of the electron must be in the opposite direction compared to that of the positron.

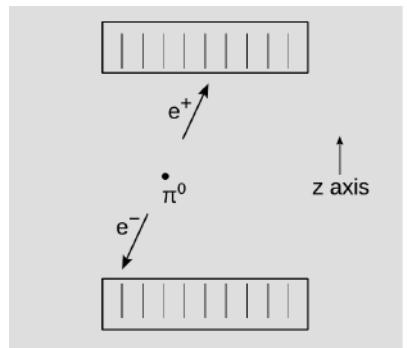
The electron and positron fly off in opposite directions due to conservation of momentum, and they could be detected by two different particle detectors lying at macroscopic distances from the place where the decay happened, as in figure b. Although separated, they are entangled. Suppose each of the detectors is capable of detecting the component of the spin along a z axis that is defined by the orientation of the detector itself. For example, the detector could in principle be a Stern-Gerlach spectrometer (sec. 14.1, p. 959), although in practice some other, more efficient method would be used. If one detector measures $s_z = +1/2$, then the other is guaranteed to see $s_z = -1/2$, because anything else would violate conservation of angular momentum. That is, the wavefunction of the system is of the form

$$\Psi = c|\uparrow\downarrow\rangle + c'|\downarrow\uparrow\rangle,$$

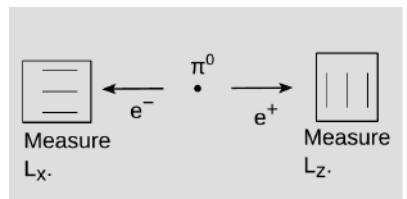
where normalization requires that $c^2 + c'^2 = 1$. If we had some way to point the pion in a certain direction before it decayed, or produce it so that it was pointed in a certain direction, then perhaps we could have arranged things so that one of the two possibilities, say $|\uparrow\downarrow\rangle$, was more likely. But the pion has spin 0, and a spinless particle is like a perfectly smooth and featureless ping-pong ball; there is no way to impose, define, or measure an orientation for it. Therefore by symmetry we have $c^2 = c'^2$. For example, we could have c and c' both equal to $1/\sqrt{2}$, or $c = i/\sqrt{2}$ and $c' = -1/\sqrt{2}$. The states $|\uparrow\downarrow\rangle$ and $|\downarrow\uparrow\rangle$ are separable in terms of the two spins, but Ψ is entangled. In the state Ψ , neither spin has a definite value, but measuring one spin determines the other spin.

In quantum computing, once a quantum computer has started running, all of its qubits will in general be entangled with one another. That means that if we read out one qubit, then later readouts of other qubits will have results that are correlated with what we got when we read out the first one. With classical information, we can always do things like splitting a book up into chapters, or distributing a long movie on two DVDs. That doesn't always work for a quantum computer. It *might* work if part of the data was separable from another part, but we would need a computer program to scan through the data and figure out whether this was in fact possible. This is called the separability problem, and unfortunately it is known to be intractable.

The no-cloning theorem described on p. 1005 is only a prohibition on making a *separable* copy of an unknown state. To see



b / The decay of a neutral pion is detected through its decay products.



c / An attempt to measure L_x and L_z simultaneously.

why, consider an experiment like the one in figure c, in which we set up the detectors so that their spin-detecting axes are in perpendicular orientations. Say one detector measures the spin of the electron along the x axis, while the other measures the positron's z spin. Now it seems that we can infer simultaneous values of both L_x and L_z for each particle, but that is impossible because L_x and L_z are incompatible observables (p. 924). Well, suppose that we measure the electron's L_x first, and then the positron's L_z . This is actually equivalent to measuring $-L_x$ for the the *positron*, and then L_z for the positron. No paradox arises, because one of the measurements will inevitably have changed the positron's spin. Going back to the version of the experiment using the entangled electron and positron, the same thing happens. For example, measuring the electron's spin has the ability to change the *positron's* spin, because they're entangled. The no-cloning theorem cannot possibly prohibit making *entangled* copies, because then it would forbid entanglement itself. Only making *separable* copies inevitably leads to paradoxes.

Problems

The symbols \checkmark , \blacksquare , etc. are explained on page 1017.

1 Nearly all naturally occurring oxygen nuclei are the isotope ^{16}O . The extremely neutron-rich isotope ^{22}O has been produced in accelerator experiments, but only with great difficulty, and little is known about its properties. The only states that have been observed and assigned reliable spins are the ground state, with spin 0, and an excited state with spin 2 and an excitation energy of 3.2 MeV. The excited state was detected by observing gamma rays for the $2 \rightarrow 0$ transition. On the hypothesis that the spin-2 excited state is a rotation, predict the gamma-ray energy that experimentalists should expect from the $4 \rightarrow 2$ transition in the same rotational band.

\checkmark \blacksquare

2 For vectors in two dimensions, which of the following are possible choices of a basis?

$$\{\hat{\mathbf{x}}\} \quad \{\hat{\mathbf{x}}, \hat{\mathbf{y}}\} \quad \{-\hat{\mathbf{x}}, \hat{\mathbf{x}} + \hat{\mathbf{y}}\} \quad \{\hat{\mathbf{x}}, \hat{\mathbf{y}}, \hat{\mathbf{x}} + \hat{\mathbf{y}}\}$$

\triangleright Solution, p. 1056 \blacksquare

3 (a) Consider the set of vectors in two dimensions. This set P is a vector space, and can be visualized as a plane, with each vector being like an arrow that extends from the origin to a particular point. Now consider the line ℓ defined by the equation $y = x$ in Cartesian coordinates, and the ray r defined by $y = x$ with $x \geq 0$. Sketch ℓ and r . If we consider ℓ and r as subsets of the arrows in P , is ℓ a vector space? Is r ?

(b) Consider the set C of angles $0 \leq \theta < 2\pi$. Define addition on C by adding the angles and then, if necessary, bringing the result back into the required range. For example, if $x = \pi$ and $y = 3\pi/2$, then $x+y = \pi/2$. Thus if we visualize C as a circle, every point on the circle has a single number to represent it, not multiple representations such as $\pi/2$ and $5\pi/2$. Suppose we want to make C into a vector space over the real numbers, so that elements of C are the vectors, while a scalar α can be *any* real number, not just a number from 0 to 2π . Then for example if $\alpha = 2$ is a scalar and $v = \pi$ is a vector, then $\alpha v = 0$. Find an example to prove that C is not a vector space, because it violates the distributive property $\alpha(v+w) = \alpha v + \alpha w$.

\triangleright Solution, p. 1057 \blacksquare

4 In the SI, we have three base units, the kilogram, the meter, and the second. From these, we form expressions such as m/s to represent units of velocity, and $\text{kg} \cdot \text{m}/\text{s}^2$ for force. Show that these expressions form a vector space with the rational numbers as the scalars. What operation on the units should we take as the “addition” operation? What operation should scalar “multiplication” be?

\triangleright Solution, p. 1057 \blacksquare

5 In problem 23 on p. 946, you showed that a wavefunction of the form

$$\Psi_0(x) = e^{-x^2/2}$$

was a solution of the Schrödinger equation for the quantum harmonic oscillator in one dimensions. (We ignore units, and the factor of 1/2 in the exponent is just a convention.) It represents the ground state. The wavefunction of the first excited state is

$$\Psi_1(x) = xe^{-x^2/2},$$

with the same value of b .

(a) Show that these states are orthogonal in the sense defined on p. 985.

(b) What is an observable that would distinguish them? ■

6 (a) When an excited state in a nucleus undergoes gamma decay, the half-life depends on a variety of factors, but a fairly typical value would be about 1 ns. Find the uncertainty in energy imposed by the energy-time uncertainty relation, and compare with a typical excitation energy of 1 MeV.

(b) Some very neutron-rich nuclei are unstable with respect to emission of a neutron, and in these cases the half-life is typically on the order of 10^{-21} s. Carry out an estimate as in part a. ■

7 As you might have guessed from the equations given in problem 5, the m th excited state of the one-dimensional quantum harmonic oscillator has a wavefunction of the form

$$\Psi_m(x) = H_m(x)e^{-x^2/2}.$$

Here H_m is a polynomial of order m , and H_m is an even function if m is even, odd if m is odd. Given these assumptions, it is possible to find Ψ_2 simply from the requirement that it be orthogonal to Ψ_0 and Ψ_1 , without having to solve the Schrödinger equation. Find H_2 by this method. (Don't worry about normalization or phase.) Hint: Near the end of the calculation, you will encounter integrals of the form $\int_{-\infty}^{\infty} x^m e^{-x^2} dx$. This can be done using software, or you can use integration by parts to relate this integral to the corresponding integral for $m - 2$. ✓ ■

8 In example 15 on p. 992, we defined a very naughty energy operator

$$\hat{H}\Psi = i\Psi.$$

Show that it is not hermitian, by directly using the definition on p. 986. ■

9 This problem refers to the analysis of the ammonia molecule in sec. 14.7.2, p. 995. (a) The bond lengths in this molecule are on the order of 0.1 nm. Use this fact to estimate the moment of inertia for rotation about the symmetry axis, and verify that states with $L_z > 0$ are likely to be populated at room temperature.

(b) The original 1955 paper by Townes and Schawlow on the microwave spectroscopy of ammonia detected about 55 lines lying between 17 and 29 GHz. Each of these corresponds to a certain value of L and L_z . Since there are many lines crowded together in this region of the spectrum, the issue arises of whether the resolution of the experiment will be sufficient to distinguish them. One of the factors limiting the resolution is that the molecules of ammonia gas have velocities that are random and randomly oriented, and this causes random Doppler shifts in the lines. Estimate the Doppler shifts at room temperature and determine whether or not they are likely to cause problems. ■

10 This problem refers to the analysis of the ammonia molecule in sec. 14.7.2, p. 995. (a) The text constructs the ground state $|g.s.\rangle$, which has energy $-|f| = f$. Use the same method to find the excited state, which has energy $+|f| = -f$.

(b) Verify that these two states are orthogonal.
(c) Find normalized versions of the two states. ■

11 Consider the wavefunctions $\Psi_1 = \nearrow$ and $\Psi_2 = \nwarrow$ for a particle in a one-dimensional box. Suppose we have the superposition $\Psi = A(2\Psi_1 + \Psi_2)$.

- (a) If Ψ is to be properly normalized, what is $|A|?$ ✓
- (b) Sketch the wavefunction.
- (c) Suppose you can measure the position of the particle very accurately. What is the probability that the particle will be found in the left half of the box? ✓
- (d) Instead of measuring position, suppose you measure the energy of the state. What is the probability that you'll measure the ground state energy? ✓
- (e) Suppose that the wavefunction had been $\Phi = A(2\Psi_1 - \Psi_2)$. Which of your answers to parts a-d would remain the same, and which would change? (You need not redo the work for the ones that would change. Just give your reasoning as to whether they would or would not.) [Problem by B. Shotwell.] ■

12 This problem builds on the results of problems 13-21 (p. 946) and 13-38 (p. 949).

Suppose we have a three-dimensional box of dimensions $L \times L \times L/2$. Let the box be oriented so that the shorter dimension is along the z direction. For convenience, define the quantity $\epsilon = h^2/8mL^2$, which has units of energy.

(a) What are the five lowest energies allowed in this box, expressed in terms of ϵ ? Give the quantum numbers for each energy, and find the degeneracy (p. 922) of each.

(b) Suppose we put five electrons in this box such that they have the lowest possible total energy. (Keep in mind that there is a limit to how many electrons can have the same spatial wavefunction.) What is the total energy of this state? \checkmark

(c) What are the two lowest-energy photons that can excite one of the five electrons (from the situation described in part b) to an excited state? \checkmark

[Problem by B. Shotwell.] \blacksquare

13 In a helium nucleus, each particle feels a potential due to the attractive forces from the other three. This potential can be well approximated as

$$U = \frac{1}{2}kr^2,$$

where k is a constant and r is the distance from the center of mass in three dimensions. You will need the result of problem 23, p. 946. Refer also to sec. 14.2.4, p. 964.

(a) Show that the Schrödinger equation is separable in terms of x , y , and z in this example.

(b) Find the ground-state wavefunction, expressed in terms of r and the constant b defined in problem 23, p. 946. \checkmark

(c) Find the energy of the ground state, expressed in terms of the classical frequency ω . \checkmark \blacksquare

14 An entangled state of three particles is prepared, described by the wavefunction,

$$\Psi = k [-2| \uparrow\uparrow\rangle + 2i| \uparrow\downarrow\rangle + | \downarrow\downarrow\rangle]$$

where the arrows are the z -components of the spins of particles A, B, and C, respectively.

(a) The wavefunction is not properly normalized. What value of $|k|$ will normalize the wavefunction? \checkmark

(b) What is the probability that measuring the z -component of particle A's spin will give spin up? \checkmark

(c) Suppose that we first measure the spin of particle C and find that it is down, and then we measure the spin of particle A. What is the probability that we will find A's spin to be up? \checkmark

[Problem by B. Shotwell.] \blacksquare

15 Suppose that we replace the usual probability rule of quantum mechanics with one of the form $P \propto |\langle \Psi | \Psi \rangle|^M$, with $M > 2$. Suppose $M = 4$. Show, by considering the example in discussion question B on p. 919, that this leads to nonconservation of probability. \blacksquare

16 In example 13 on p. 988, we defined unnormalized wavefunctions for the traveling-wave solutions to the Schrödinger equation in a “quantum moat,” and calculated the inner product $\langle \text{ccw} | \text{cw} \rangle = 0$ to verify that the counterclockwise and clockwise traveling waves were orthogonal, as must be the case for distinguishable states. Suppose we want to define a normalized version of the counterclockwise wave, $|\text{ccw}\rangle = Ae^{i\theta}$. Use an inner product to determine $|A|^2$, and show that normalization doesn’t depend on the phase of A . (Do not assume that A is real.) $\checkmark \blacksquare$

17 In example 13 on p. 988, we defined wavefunctions for the traveling-wave solutions to the Schrödinger equation in a “quantum moat,” and calculated the inner product $\langle \text{ccw} | \text{cw} \rangle = 0$ to verify that the counterclockwise and clockwise traveling waves were orthogonal, as must be the case for distinguishable states. Let’s now define standing-wave versions $|c\rangle = \cos \theta$ and $|s\rangle = \sin \theta$. Verify by direct calculation that $\langle c | s \rangle = 0$.

Remark: Note that, as discussed in the sidebar on p. 969, this does not contradict the principle that a quantum-mechanical phase is undetectable. \blacksquare

18 Consider the wavefunctions

$$\begin{aligned}\Psi_1 &= e^{ikx}, \\ \Psi_2 &= e^{-ikx}, \\ \Phi_1 &= \cos kx, \quad \text{and} \\ \Phi_2 &= i \sin kx.\end{aligned}$$

Show that $\Psi_1 = \Phi_1 + \Phi_2$. Similarly, express Ψ_2 in terms of the Φ ’s, and express each of the Φ ’s in terms of the Ψ ’s. Relate this to the principle that there is no preferred basis in quantum mechanics (p. 990). $\checkmark \blacksquare$

19 In section 14.5.1, p. 973, we found the solution to the Schrödinger equation for a particle arriving at a potential barrier, in the case where the far side of the barrier is classically allowed. We now consider the case in which the barrier is not just high enough to make the region beyond it classically forbidden — we let the height of the barrier be infinite. Let the potential be

$$U(x) = \begin{cases} 0, & x < 0 \\ +\infty, & x > 0, \end{cases}$$

and let the incident wave be

$$\Psi_I = e^{i(kx - \omega t)} \quad (x < 0).$$

Determine the form of the complete solution to the Schrödinger equation on the left side of the barrier, including both the incident wave and the reflected wave. \checkmark \blacksquare

20 As discussed on p. 989, suppose that for a particle in a box, we have

$$\begin{aligned} \mathcal{O}_E \curvearrowleft &= 1 \curvearrowleft, \\ \mathcal{O}_E \curvearrowright &= 4 \curvearrowright, \\ \Psi &= c \curvearrowleft + c' \curvearrowright, \quad \text{and} \\ |c| &= |c'|. \end{aligned}$$

Show that $\langle \Psi | \mathcal{O}_E \Psi \rangle = 2.5$. \triangleright Solution, p. 1057 \blacksquare

21 Consider the wavefunctions

$$\begin{aligned} \Psi_1 &= e^{x+it} \quad \text{and} \\ \Psi_2 &= e^{t+ix}. \end{aligned}$$

To keep the writing simple, we use a system of units (not SI) such that these expressions make sense, and in which $\hbar = 1$.

- (a) Show by direct substitution in the time-dependent Schrödinger equation (with $U = \text{constant}$) that one of these is a solution and the other is not.
- (b) Make an independent argument, requiring no calculations, to the effect that the invalid one violates one of the fundamental principles 1-5 of quantum mechanics listed in section 14.6.5, p. 993. \blacksquare

22 Microscopic circuits are etched on the surface of a silicon chip. The equivalent of a wire in such an integrated circuit is called a “trace.” We consider the case where the trace is narrow enough to make quantum effects relevant, and we treat an electron inside the trace using the two-dimensional Schrödinger equation. We describe the trace as an infinite strip running parallel to the x axis, extending from $y = 0$ to $y = b$. The potential is

$$U = \begin{cases} 0, & 0 < y < b \\ +\infty, & y \leq 0 \text{ or } y \geq b \end{cases}$$

For convenience of notation, let $a = \pi/b$. Consider the following wavefunctions:

$$\begin{aligned}\Psi_1 &= e^{i(-kx-\omega t)} \sin ay \\ \Psi_2 &= e^{-i\omega t} \sin kx e^{ay} \\ \Psi_3 &= e^{-i\omega t} e^{kx} \sin ay \\ \Psi_4 &= e^{-i\omega t} (\sin 2ax \sin ay + \sin ax \sin 2ay)\end{aligned}$$

The symbols ω and k stand for real constants. Identify the wavefunctions that have the following properties. Exactly one of the wavefunctions has each property. Explain all answers.

- (a) cannot be a solution of the Schrödinger equation for this potential
- (b) is a traveling wave solution
- (c) is a solution that could represent the case where $0 < y < b$ is classically forbidden
- (d) is not separable



Key to symbols:

■ easy ■ typical ■ challenging ■ difficult ■ very difficult

✓ An answer check is available at www.lightandmatter.com.

Three essential mathematical skills

More often than not when a search-and-rescue team finds a hiker dead in the wilderness, it turns out that the person won a Darwin Award by not carrying some item from a short list of essentials, such as water and a map. There are three mathematical essentials in this course.

1. Converting units

basic technique: subsection 0.1.9, p. 28; conversion of area, volume, etc.: subsection 0.2.1, p. 34

Examples:

$$0.7 \cancel{\text{kg}} \times \frac{10^3 \text{ g}}{1 \cancel{\text{kg}}} = 700 \text{ g}.$$

To check that we have the conversion factor the right way up (10^3 rather than $1/10^3$), we note that the smaller unit of grams has been *compensated* for by making the number larger.

For units like m^2 , kg/m^3 , etc., we have to raise the conversion factor to the appropriate power:

$$4 \text{ m}^3 \times \left(\frac{10^3 \text{ mm}}{1 \text{ m}} \right)^3 = 4 \times 10^9 \cancel{\text{m}^3} \times \frac{\text{mm}^3}{\cancel{\text{m}^3}} = 4 \times 10^9 \text{ mm}^3$$

Examples with solutions — p. 47, #6; p. 51, #31

Problems you can check at lightandmatter.com/area1checker.html — p. 47, #5; p. 47, #4; p. 47, #7; p. 51, #22; p. 52, #40

2. Reasoning about ratios and proportionalities

The technique is introduced in subsection 0.2.2, p. 35, in the context of area and volume, but it applies more generally to any relationship in which one variable depends on another raised to some power.

Example: When a car or truck travels over a road, there is wear and tear on the road surface, which incurs a cost. Studies show that the cost per kilometer of travel C is given by

$$C = kw^4,$$

where w is the weight per axle and k is a constant. The weight per axle is about 13 times higher for a semi-trailer than for my Honda Fit. How many times greater is the cost imposed on the federal government when the semi travels a given distance on an interstate freeway?

▷ First we convert the equation into a proportionality by throwing out k , which is the same for both vehicles:

$$C \propto w^4$$

Next we convert this proportionality to a statement about ratios:

$$\frac{C_1}{C_2} = \left(\frac{w_1}{w_2} \right)^4 \approx 29,000$$

Since the gas taxes paid by the trucker are nowhere near 29,000 times more than those I pay to drive my Fit the same distance, the federal government is effectively awarding a massive subsidy to the trucking company. Plus my Fit is cuter.

Examples with solutions — p. 51, #32; p. 51, #33; p. 52, #38; p. 49, #17

Problems you can check at lightandmatter.com/area1checker.html — p. 52, #37; p. 52, #39; p. 123, #24; p. 124, #27; p. 121, #9; p. 294, #3

3. Vector addition

subsection 3.4.3, p. 203

Example: The $\Delta\mathbf{r}$ vector from San Diego to Los Angeles has magnitude 190 km and direction 129° counterclockwise from east. The one from LA to Las Vegas is 370 km at 38° counterclockwise from east. Find the distance and direction from San Diego to Las Vegas.

▷ Graphical addition is discussed on p. 203. Here we concentrate on analytic addition, which involves adding the x components to find the total x component, and similarly for y . The trig needed in order to find the components of the second leg (LA to Vegas) is laid out in figure 1 on p. 201 and explained in detail in example 60 on p. 201:

$$\Delta x_2 = (370 \text{ km}) \cos 38^\circ = 292 \text{ km}$$

$$\Delta y_2 = (370 \text{ km}) \sin 38^\circ = 228 \text{ km}$$

(Since these are intermediate results, we keep an extra sig fig to avoid accumulating too much rounding error.) Once we understand the trig for one example, we don't need to reinvent the wheel every time. The pattern is completely universal, provided that we first make sure to get the angle expressed according to the usual trig convention, counterclockwise from the x axis. Applying the pattern to the first leg, we have:

$$\Delta x_1 = (190 \text{ km}) \cos 129^\circ = -120 \text{ km}$$

$$\Delta y_1 = (190 \text{ km}) \sin 129^\circ = 148 \text{ km}$$

For the vector directly from San Diego to Las Vegas, we have

$$\Delta x = \Delta x_1 + \Delta x_2 = 172 \text{ km}$$

$$\Delta y = \Delta y_1 + \Delta y_2 = 376 \text{ km.}$$

The distance from San Diego to Las Vegas is found using the Pythagorean theorem,

$$\sqrt{(172 \text{ km})^2 + (376 \text{ km})^2} = 410 \text{ km}$$

(rounded to two sig figs because it's one of our final results). The direction is one of the two possible values of the inverse tangent

$$\tan^{-1}(\Delta y/\Delta x) = \{65^\circ, 245^\circ\}.$$

Consulting a sketch shows that the first of these values is the correct one.

Examples with solutions — p. 233, #62; p. 233, #63; p. 230, #45

Problems you can check at lightandmatter.com/area1checker.html — p. 232, #53; p. 232, #57; p. 233, #58; p. 238, #79; p. 230, #46

Mathematical Review

Algebra

Quadratic equation:

The solutions of $ax^2 + bx + c = 0$
are $x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$.

Logarithms and exponentials:

$$\ln(ab) = \ln a + \ln b$$

$$e^{a+b} = e^a e^b$$

$$\ln e^x = e^{\ln x} = x$$

$$\ln(a^b) = b \ln a$$

$$\frac{d}{dx}(cf) = c \frac{df}{dx}$$

$$\frac{d}{dx}(f + g) = \frac{df}{dx} + \frac{dg}{dx}$$

The chain rule:

$$\frac{d}{dx}f(g(x)) = f'(g(x))g'(x)$$

Geometry, area, and volume

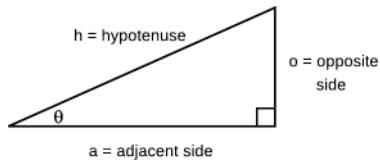
area of a triangle of base b and height h
circumference of a circle of radius r
area of a circle of radius r
surface area of a sphere of radius r
volume of a sphere of radius r

$$\begin{aligned} &= \frac{1}{2}bh \\ &= 2\pi r \\ &= \pi r^2 \\ &= 4\pi r^2 \\ &= \frac{4}{3}\pi r^3 \end{aligned}$$

$$\frac{d}{dx}(fg) = \frac{df}{dx}g + \frac{dg}{dx}f$$

$$\frac{d}{dx}\left(\frac{f}{g}\right) = \frac{f'}{g} - \frac{fg'}{g^2}$$

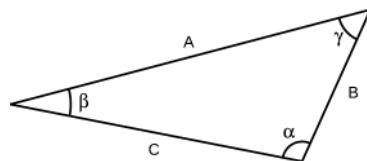
Trigonometry with a right triangle



$$\sin \theta = o/h \quad \cos \theta = a/h \quad \tan \theta = o/a$$

$$\text{Pythagorean theorem: } h^2 = a^2 + o^2$$

Trigonometry with any triangle



Law of Sines:

$$\frac{\sin \alpha}{A} = \frac{\sin \beta}{B} = \frac{\sin \gamma}{C}$$

Law of Cosines:

$$C^2 = A^2 + B^2 - 2AB \cos \gamma$$

Properties of the derivative and integral

Let f and g be functions of x , and let c be a constant.

Linearity of the derivative:

$$\frac{d}{dx}(cf) = c \frac{df}{dx}$$

$$\frac{d}{dx}(f + g) = \frac{df}{dx} + \frac{dg}{dx}$$

The chain rule:

$$\frac{d}{dx}f(g(x)) = f'(g(x))g'(x)$$

Derivatives of products and quotients:

Some derivatives:

$$\begin{aligned} \frac{d}{dx}x^m &= mx^{m-1}, \text{ except for } m = 0 \\ \frac{d}{dx}\sin x &= \cos x & \frac{d}{dx}\cos x &= -\sin x \\ \frac{d}{dx}e^x &= e^x & \frac{d}{dx}\ln x &= \frac{1}{x} \end{aligned}$$

Linearity of the integral:

$$\int cf(x) dx = c \int f(x) dx$$

$$\int [f(x) + g(x)] dx = \int f(x) dx + \int g(x) dx$$

The fundamental theorem of calculus:

The derivative and the integral undo each other, in the following sense:

$$\int_a^b f'(x) dx = f(b) - f(a)$$

Approximations to Exponents and Logarithms

It is often useful to have certain approximations involving exponents and logarithms. As a simple numerical example, suppose that your bank balance grows by 1% for two years in a row. Then the result of compound interest is growth by a factor of $1.01^2 = 1.0201$, but the compounding effect is quite small, and the result is essentially 2% growth. That is, $1.01^2 \approx 1.02$. This is a special case of the more general approximation

$$(1 + \epsilon)^p \approx 1 + p\epsilon,$$

which holds for small values of ϵ and is used in example 4 on p. 408 relating to relativity. Proof: Any real exponent p can be approximated to the desired precision as $p = a/b$, where a and b are integers. Let $(1 + \epsilon)^p = 1 + x$. Then $(1 + \epsilon)^a = (1 + x)^b$. Multiplying out both sides gives $1 + a\epsilon + \dots = 1 + bx + \dots$, where \dots indicates higher powers. Neglecting these higher powers gives $x \approx (a/b)\epsilon \approx p\epsilon$.

We have considered an approximation that can be found by restricting the *base* of an exponential to be close to 1. It is often of interest as well to consider the case where the *exponent* is restricted to be small. Consider the base- e case. One way of defining e is that when we use it as a base, the rate of growth of the function e^x , for small x , equals 1. That is,

$$e^x \approx 1 + x$$

for small x . This can easily be generalized to other bases, since $a^x = e^{\ln(a^x)} = e^{x \ln a}$, giving

$$a^x \approx 1 + x \ln a.$$

Finally, since $e^x \approx 1 + x$, we also have

$$\ln(1 + x) \approx x.$$

Programming with python

The purpose of this tutorial is to help you get familiar with a computer programming language called Python, which I've chosen because (a) it's free, and (b) it's easy to use interactively. I won't assume you have any previous experience with computer programming; you won't need to learn very much Python, and what little you do need to learn I'll explain explicitly. If you really want to learn Python more thoroughly, there are a couple of excellent books that you can download for free on the Web:

How to Think Like a Computer Scientist (Python Version), Allen B. Downey, Jeffrey Elkner, Moshe Zadka,
<http://www.ibiblio.org/obp/>

Dive Into Python, Mark Pilgrim,
<http://diveintopython.net/>

The first book is meant for people who have never programmed before, while the second is a more complete introduction aimed at veteran programmers who know a different language already.

Using Python as a calculator

The easiest way to get Python going is to go to the web site ideone.com. Under "choose a language," select Python. Inside the window where it says "paste your source code or insert template or sample," type `print(2+2)`. Click on the "submit" button. The result, 4, is shown under "output." In other words, you can use Python just like a calculator.

For compactness, I'll show examples in the following style:

```
>>> print(2+2)
4
```

Here the `>>>` is not something you would type yourself; it's just a marker to distinguish your input from the program's output. (In some

older versions of Python, the computer will actually print out `>>>` as a prompt to tell you it's ready to type something.)

There are only a couple of things to watch out for. First, Python distinguishes between integers and real numbers, so the following gives an unexpected result:

```
>>> print(2/3)
0
```

To get it to treat these values as real numbers, you have to use decimal points:

```
>>> print(2./3.)
0.6666666666666666666666666663
```

Multiplication is represented by `*`:

```
>>> print(2.*3.)
6.0
```

Also, Python doesn't know about its own library of math functions unless you tell it explicitly to load them in:

```
>>> print (sqrt(2.))
Traceback (most recent call last):
File '<stdin>', line 1, in ?
NameError: There is no variable named 'sqrt'
```

Here are the steps you have to go through to calculate the square root of 2 successfully:

```
>>> import math
>>> print(math.sqrt(2.))
1.4142135623730951
```

The first line is just something you can make a habit of doing every time you start up Python. In the second line, the name of the square root function had to be prefixed with `"math."` to tell Python where you wanted to get this `"sqrt"` function from. (All of this may seem like a nuisance if you're just using Python as a

calculator, but it's a good way to design a programming language so that names of functions never conflict.)

Try it. Experiment and figure out whether Python's trig functions assume radians or degrees.

Variables

Python lets you define variables and assign values to them using an equals sign:

```
>>> dwarfs=7
>>> print(dwarfs)
>>> print(dwarfs+3)
7
10
```

Note that a variable in computer programming isn't quite like a variable in algebra. In algebra, if $a=7$ then $a=7$ always, throughout a particular calculation. But in a programming language, the variable name really represents a place in memory where a number can be stored, so you can change its value:

```
>>> dwarfs=7
>>> dwarfs=37
>>> print(dwarfs)
37
```

You can even do stuff like this,

```
>>> dwarfs=37
>>> dwarfs=dwarfs+1
>>> print(dwarfs)
38
```

In algebra it would be nonsense to have a variable equal to itself plus one, but in a computer program, it's not an assertion that the two things are equal, its a command to calculate the value of the expression on the right side of the equals, and then put that number into the memory location referred to by the variable name on the left.

Try it. What happens if you do $dwarfs+1 = dwarfs$? Do you understand why?

Functions

Somebody had to teach Python how to do functions like `sqrt`, and it's handy to be able to define your own functions in the same way. Here's how to do it:

```
>>> def double(x):
>>>     return 2.*x
>>> print(double(5.))
10.0
```

Note that the indentation is mandatory. The first and second lines define a function called `double`. The final line evaluates that function with an input of 5.

Loops

Suppose we want to add up all the numbers from 0 to 99.

Automating this kind of thing is exactly what computers are best at, and Python provides a mechanism for this called a loop:

```
>>> sum=0
>>> for j in range(100):
>>>     sum=sum+j
>>> print(sum)
4950
```

The stuff that gets repeated — the inside of the loop — has to be indented, just like in a function definition. Python always counts loops starting from 0, so for `j in range(100)` actually causes `j` to range from 0 to 99, not from 1 to 100.

Appendix 2: Miscellany

Unphysical “hovering” solutions to conservation of energy

On page 83, I gave the following derivation for the acceleration of an object under the influence of gravity:

$$\begin{aligned}\left(\frac{dv}{dt}\right) &= \left(\frac{dv}{dK}\right) \left(\frac{dK}{dU}\right) \left(\frac{dU}{dy}\right) \left(\frac{dy}{dt}\right) \\ &= \left(\frac{1}{mv}\right) (-1)(mg)(v) \\ &= -g\end{aligned}$$

There is a loophole in this argument, however. When I say $dv/dK = 1/(mv)$, that only works when the object is moving. If it’s at rest, v is nondifferentiable as a function of K (or we could say that the derivative is infinite). Energy can in fact be conserved by an object that simply hovers above the ground: its kinetic energy is constant, and its gravitational energy is also constant. Why, then, do we never observe such behavior, except in Coyote and Roadrunner cartoons when the Coyote runs off the edge of a cliff without noticing it at first?

Suppose we toss a baseball straight up, and pick a coordinate system in which upward velocities are positive. The ball’s velocity is a continuous function of time, and it changes from being positive to being negative, so there must be some instant at which it equals zero. Conservation of energy would be satisfied if the velocity were to remain at zero for a minute or an hour before the ball finally made the decision to fall. One thing that seems odd about all this is that there’s no obvious way for the ball to “decide” when it was time to go ahead and fall back down again. It violates the principle that the laws of physics are supposed to be deterministic.

One reason that we could never hope to observe such behavior in reality is that the ball would have to spend some time being *exactly* at rest, and yet no object can ever stay exactly at rest for any finite amount of time. Objects in the real world are buffeted by air currents, for example. At the atomic level, the interaction of these air currents with the ball consists of discrete collisions with whizzing air molecules, and a quick back-of-the-envelope estimate shows that for an object this size, the typical time between collisions is on the order of 10^{-27} s, which would limit the duration of the hovering to a time far too short to allow it to be observed. Nevertheless this is not a completely satisfying explanation. It makes us wonder whether we ought to apply to the government for a research grant to do an experiment in which a baseball would be shot upward in a chamber that had been pumped out to an ultra-high vacuum!

A somewhat better approach is to consider that motion is relative, so the ball’s velocity can only be zero in one particular frame of reference. It wouldn’t make sense for the ball to exhibit qualitatively different behavior when it was at rest, because different observers don’t even agree when the ball is at rest. But this argument also fails to resolve the issue completely, because this is a ball interacting with the planet Earth via gravitational forces, so it could make a difference whether the ball was at rest *relative to the earth*. Suppose we go into a frame of reference defined by an observer watching the ball as she descends in a glass elevator. At the moment when the

ball is at rest relative to the earth, she sees both the ball and the earth as moving upward at the same speed. It would be perfectly consistent with conservation of energy if she were to see them maintain this distance from one another for several minutes. In her frame, their kinetic energies would be nonzero, but constant, and the gravitational energy only depends on the separation between the ball and the earth, so it would be constant as well.

Now that we're thinking of the ball and the earth as two objects interacting with one another, it becomes natural to think of them on the same footing. What about the motion of the Earth? The earth feels a gravitational attraction from the ball, just as the ball feels one from the Earth. To make this symmetry more evident, let's imagine two planets of equal mass, Foo and Bar, initially at rest with respect to one another. The Fooites and Barians realize that the gravitational interaction between their planets will cause them to drop together and collide. It seems that they should get ready for the end of the world. And yet before they riot, get drunk, or tell their spouses that yes, they really *do* look fat in that dress, maybe they should consider the possibility that the two planets will simply hover in place for some amount of time, because that would satisfy conservation of energy. Now the physical implausibility of the hovering solution becomes even more apparent. Not only does one planet have to "decide" at precisely what microsecond to go ahead and fall, but the other planet has to make the same decision at the same instant, or else conservation of energy will be violated. There is no physical process or interaction between the two planets that could perfectly synchronize their "decisions" like this. (The mechanism can't be gravity, because nothing about the gravitational interaction provides any kind of a count-down that would pick out one particular time as the one at which the planets should start moving.)

The key to making sense of all this is to realize that each planet can only "feel" the gravitational field in its own region of space. Its acceleration can only depend on the field, and not on the detailed arrangement of masses elsewhere in the universe that caused that field. Granting this kind of "real" status to fields can be considered as a logically necessary supplement to conservation of energy.

Automated search for the brachistochrone

See page 95.

```
1  d=.01
2  c1=.61905
3  c2=-.94427
4  a = 1.
5  b = 1.
6  for i in range(100):
7      bestt = 99.
8      for j in range(3):
9          for k in range(3):
10             try_c1 = c1+(j-1)*d
11             try_c2 = c2+(k-1)*d
12             t = timeb(a,b,try_c1,try_c2,100000)
13             if t<bestt :
14                 bestc1 = try_c1
15                 bestc2 = try_c2
16                 bestj = j
17                 bestk = k
```

```

18      bestt = t
19      c1 = bestc1
20      c2 = bestc2
21      c3 = (b-c1*a-c2*a**2)/(a**3)
22      print(c1, c2, c3, bestt)
23      if (bestj == 1) and (bestk == 1) :
24          d = d*.5

```

Derivation of the steady state for damped, driven oscillations

Using the trig identities for the sine of a sum and cosine of a sum, we can change equation [2] on page 181 into the form

$$\begin{aligned} & [(-m\omega^2 + k) \cos \delta - b\omega \sin \delta - F_m/A] \sin \omega t \\ & + [(-m\omega^2 + k) \sin \delta + b\omega \cos \delta] \cos \omega t = 0. \end{aligned}$$

Both the quantities in square brackets must equal zero, which gives us two equations we can use to determine the unknowns A and δ . The results are

$$\begin{aligned} \delta &= \tan^{-1} \frac{b\omega}{m\omega^2 - k} \\ &= \tan^{-1} \frac{\omega\omega_0}{Q(\omega_0^2 - \omega^2)} \end{aligned}$$

and

$$\begin{aligned} A &= \frac{F_m}{\sqrt{(m\omega^2 - k)^2 + b^2\omega^2}} \\ &= \frac{F_m}{m\sqrt{(\omega^2 - \omega_0^2)^2 + \omega_0^2\omega^2 Q^{-2}}}. \end{aligned}$$

Proofs relating to angular momentum

Uniqueness of the cross product

The vector cross product as we have defined it has the following properties:

- (1) It does not violate rotational invariance.
- (2) It has the property $\mathbf{A} \times (\mathbf{B} + \mathbf{C}) = \mathbf{A} \times \mathbf{B} + \mathbf{A} \times \mathbf{C}$.
- (3) It has the property $\mathbf{A} \times (k\mathbf{B}) = k(\mathbf{A} \times \mathbf{B})$, where k is a scalar.

Theorem: The definition we have given is the only possible method of multiplying two vectors to make a third vector which has these properties, with the exception of trivial redefinitions which just involve multiplying all the results by the same constant or swapping the names of the axes. (Specifically, using a left-hand rule rather than a right-hand rule corresponds to multiplying all the results by -1 .)

Proof: We prove only the uniqueness of the definition, without explicitly proving that it has properties (1) through (3).

Using properties (2) and (3), we can break down any vector multiplication $(A_x\hat{\mathbf{x}} + A_y\hat{\mathbf{y}} + A_z\hat{\mathbf{z}}) \times (B_x\hat{\mathbf{x}} + B_y\hat{\mathbf{y}} + B_z\hat{\mathbf{z}})$ into terms involving cross products of unit vectors.

A “self-term” like $\hat{\mathbf{x}} \times \hat{\mathbf{x}}$ must either be zero or lie along the x axis, since any other direction would violate property (1). If was not zero, then $(-\hat{\mathbf{x}}) \times (-\hat{\mathbf{x}})$ would have to lie in the opposite

direction to avoid breaking rotational invariance, but property (3) says that $(-\hat{\mathbf{x}}) \times (-\hat{\mathbf{x}})$ is the same as $\hat{\mathbf{x}} \times \hat{\mathbf{x}}$, which is a contradiction. Therefore the self-terms must be zero.

An “other-term” like $\hat{\mathbf{x}} \times \hat{\mathbf{y}}$ could conceivably have components in the x - y plane and along the z axis. If there was a nonzero component in the x - y plane, symmetry would require that it lie along the diagonal between the x and y axes, and similarly the in-the-plane component of $(-\hat{\mathbf{x}}) \times \hat{\mathbf{y}}$ would have to be along the other diagonal in the x - y plane. Property (3), however, requires that $(-\hat{\mathbf{x}}) \times \hat{\mathbf{y}}$ equal $-(\hat{\mathbf{x}} \times \hat{\mathbf{y}})$, which would be along the original diagonal. The only way it can lie along both diagonals is if it is zero.

We now know that $\hat{\mathbf{x}} \times \hat{\mathbf{y}}$ must lie along the z axis. Since we are not interested in trivial differences in definitions, we can fix $\hat{\mathbf{x}} \times \hat{\mathbf{y}} = \hat{\mathbf{z}}$, ignoring peurile possibilities such as $\hat{\mathbf{x}} \times \hat{\mathbf{y}} = 7\hat{\mathbf{z}}$ or the left-handed definition $\hat{\mathbf{x}} \times \hat{\mathbf{y}} = -\hat{\mathbf{z}}$. Given $\hat{\mathbf{x}} \times \hat{\mathbf{y}} = \hat{\mathbf{z}}$, the symmetry of space requires that similar relations hold for $\hat{\mathbf{y}} \times \hat{\mathbf{z}}$ and $\hat{\mathbf{z}} \times \hat{\mathbf{x}}$, with at most a difference in sign. A difference in sign could always be eliminated by swapping the names of some of the axes, so ignoring possible trivial differences in definitions we can assume that the cyclically related set of relations $\hat{\mathbf{x}} \times \hat{\mathbf{y}} = \hat{\mathbf{z}}$, $\hat{\mathbf{y}} \times \hat{\mathbf{z}} = \hat{\mathbf{x}}$, and $\hat{\mathbf{z}} \times \hat{\mathbf{x}} = \hat{\mathbf{y}}$ holds. Since the arbitrary cross-product with which we started can be broken down into these simpler ones, the cross product is uniquely defined.

The choice of axis theorem

Theorem: Suppose a closed system of material particles conserves angular momentum in one frame of reference, with the axis taken to be at the origin. Then conservation of angular momentum is unaffected if the origin is relocated or if we change to a frame of reference that is in constant-velocity motion with respect to the first one. The theorem also holds in the case where the system is not closed, but the total external force is zero.

Proof: In the original frame of reference, angular momentum is conserved, so we have $d\mathbf{L}/dt=0$. From example 28 on page 290, this derivative can be rewritten as

$$\frac{d\mathbf{L}}{dt} = \sum_i \mathbf{r}_i \times \mathbf{F}_i,$$

where \mathbf{F}_i is the total force acting on particle i . In other words, we’re adding up all the torques on all the particles.

By changing to the new frame of reference, we have changed the position vector of each particle according to $\mathbf{r}_i \rightarrow \mathbf{r}_i + \mathbf{k} - \mathbf{u}t$, where \mathbf{k} is a constant vector that indicates the relative position of the new origin at $t = 0$, and \mathbf{u} is the velocity of the new frame with respect to the old one. The forces are all the same in the new frame of reference, however. In the new frame, the rate of change of the angular momentum is

$$\begin{aligned} \frac{d\mathbf{L}}{dt} &= \sum_i (\mathbf{r}_i + \mathbf{k} - \mathbf{u}t) \times \mathbf{F}_i \\ &= \sum_i \mathbf{r}_i \times \mathbf{F}_i + (\mathbf{k} - \mathbf{u}t) \times \sum_i \mathbf{F}_i. \end{aligned}$$

The first term is the expression for the rate of change of the angular momentum in the original frame of reference, which is zero by assumption. The second term vanishes by Newton’s third law; since the system is closed, every force \mathbf{F}_i cancels with some force \mathbf{F}_j . (If external forces act, but they add up to zero, then the sum can be broken up into a sum of internal forces and a sum of external forces, each of which is zero.) The rate of change of the angular momentum is therefore zero in the new frame of reference.

The spin theorem

Theorem: An object's angular momentum with respect to some outside axis A can be found by adding up two parts:

- (1) The first part is the object's angular momentum found by using its own center of mass as the axis, i.e. the angular momentum the object has because it is spinning.
- (2) The other part equals the angular momentum that the object would have with respect to the axis A if it had all its mass concentrated at and moving with its center of mass.

Proof: Let \mathbf{r}_{cm} be the position of the center of mass. The total angular momentum is

$$\mathbf{L} = \sum_i \mathbf{r}_i \times \mathbf{p}_i,$$

which can be rewritten as

$$\mathbf{L} = \sum_i (\mathbf{r}_{cm} + \mathbf{r}_i - \mathbf{r}_{cm}) \times \mathbf{p}_i,$$

where $\mathbf{r}_i - \mathbf{r}_{cm}$ is particle i 's position relative to the center of mass. We then have

$$\begin{aligned} \mathbf{L} &= \mathbf{r}_{cm} \times \sum_i \mathbf{p}_i + \sum_i (\mathbf{r}_i - \mathbf{r}_{cm}) \times \mathbf{p}_i \\ &= \mathbf{r}_{cm} \times \mathbf{p}_{total} + \sum_i (\mathbf{r}_i - \mathbf{r}_{cm}) \times \mathbf{p}_i \\ &= \mathbf{r}_{cm} \times m_{total} \mathbf{v}_{cm} + \sum_i (\mathbf{r}_i - \mathbf{r}_{cm}) \times \mathbf{p}_i. \end{aligned}$$

The first and second terms in this expression correspond to the quantities (2) and (1), respectively.

Different Forms of Maxwell's Equations

First we reproduce Maxwell's equations as stated on page 724, in integral form, using the SI (meter-kilogram-second) system of units, with the coupling constants written in terms of k and c :

$$\begin{aligned} \Phi_E &= 4\pi k q_{in} \\ \Phi_B &= 0 \\ \Gamma_E &= -\frac{\partial \Phi_B}{\partial t} \\ c^2 \Gamma_B &= \frac{\partial \Phi_E}{\partial t} + 4\pi k I_{through} \end{aligned}$$

Homework problem 39 on page 755 deals with rewriting these in terms of $\epsilon_0 = 1/4\pi k$ and $\mu_0 = 4\pi k/c^2$ rather than k and c .

For the reader who has been studying the optional sections giving Maxwell's equations in differential form, here is a summary:

$$\operatorname{div} \mathbf{E} = 4\pi k\rho$$

$$\operatorname{div} \mathbf{B} = 0$$

$$\operatorname{curl} \mathbf{E} = -\frac{\partial \mathbf{B}}{\partial t}$$

$$c^2 \operatorname{curl} \mathbf{B} = \frac{\partial \mathbf{E}}{\partial t} + 4\pi k \mathbf{j}$$

Although all engineering and most scientific work these days is done in the SI (mks) system, one may still encounter the older cgs (centimeter-gram-second) system, especially in astronomy and particle physics. The mechanical units in this system include the dyne ($\text{g}\cdot\text{cm}/\text{s}^2$) for force, and the erg ($\text{g}\cdot\text{cm}^2/\text{s}^2$) for energy. The system is extended to electrical units by taking $k = 1$ as a matter of definition, so the Coulomb force law is $F = q_1 q_2 / r^2$. This equation indirectly defines a unit of charge called the electrostatic unit, with $1 \text{ C} = 2.998 \times 10^9 \text{ esu}$, the factor of 2.998 arising from the speed of light. The unit of voltage is the statvolt, $1 \text{ statvolt} = 299.8 \text{ V}$. In this system, the electric and magnetic fields have the same units, dynes/esu, but to avoid confusion, magnetic fields are normally written using the equivalent unit of gauss, $1 \text{ gauss} = 1 \text{ dyne}/\text{esu} = 10^{-4} \text{ T}$. The force on a charged particle is $\mathbf{F} = q\mathbf{E} + q\frac{\mathbf{v}}{c} \times \mathbf{B}$, which differs from the mks version by the $1/c$ factor in the magnetic term. Maxwell's equations are:

$$\Phi_E = 4\pi q_{in}$$

$$\Phi_B = 0$$

$$\Gamma_E = -\frac{1}{c} \frac{\partial \Phi_B}{\partial t}$$

$$\Gamma_B = \frac{1}{c} \frac{\partial \Phi_E}{\partial t} + \frac{4\pi}{c} I_{through}$$

Appendix 3: Photo Credits

Except as specifically noted below or in a parenthetical credit in the caption of a figure, all the illustrations in this book are by under my own copyright, and are copyleft licensed under the same license as the rest of the book.

In some cases it's clear from the date that the figure is public domain, but I don't know the name of the artist or photographer; I would be grateful to anyone who could help me to give proper credit. I have assumed that images that come from U.S. government web pages are copyright-free, since products of federal agencies fall into the public domain. When "PSSC Physics" is given as a credit, it indicates that the figure is from the second edition of the textbook entitled Physics, by the Physical Science Study Committee; these are used according to a blanket permission given in the later PSSC College Physics edition, which states on the copyright page that "The materials taken from the original and second editions and the Advanced Topics of PSSC PHYSICS included in this text will be available to all publishers for use in English after December 31, 1970, and in translations after December 31, 1975."

In a few cases, I have made use of images under the fair use doctrine. However, I am not a lawyer, and the laws on fair use are vague, so you should not assume that it's legal for you to use these images. In particular, fair use law may give you less leeway than it gives me, because I'm using the images for educational purposes, and giving the book away for free. Likewise, if the photo credit says "courtesy of ...," that means the copyright owner gave me permission to use it, but that doesn't mean you have permission to use it.

Cover Eclipse: Luc Viatour, CC-BY-SA licensed.

?? Wicked witch: W.W. Denslow, 1900. Quote from The Wizard of Oz, by L. Frank Baum, 1900. **13 Mars Climate Orbiter:** NASA/JPL/CIT. **25 Standard kilogram:** Bo Bengtsen, GFDL licensed. Further retouching by Wikipedia user Greg L and by B. Crowell. **44 Jar of jellybeans:** Flickr user cbgrfx123, CC-BY-SA licensed. **45 Amphicoelias:** Wikimedia commons users Dinoguy2, Niczar, ArthurWeasley, Steveoc 86, Dropzink, and Piotr Jaworski, CC-BY-SA licensed. **55 Galaxies:** Hubble Space Telescope. Hubble material is copyright-free and may be freely used as in the public domain without fee, on the condition that NASA and ESA is credited as the source of the material. The material was created for NASA by STScI under Contract NAS5-26555 and for ESA by the Hubble European Space Agency Information Centre.. **56 Portrait of Monsieur Lavoisier and His Wife:** Jacques-Louis David, 1788. **57 Astronaut:** NASA. **62 Galileo:** Justus Sustermans, 1636. **65 Foucault and pendulum:** contemporary, ca. 1851. **65 Galileo's trial:** Cristiano Banti (1857). **71 Wind tunnel:** NASA. **73 Portrait of James Joule:** contemporary. **74 Infrared photographs:** Courtesy of M. Vollmer and K.P. M'ollmann, Univ. Appl. Sciences, Brandenburg, Germany. **81 Skateboarder:** Courtesy of J.D. Rogge. **87 Funicular railroad:** Historic American Buildings Survey, public domain. **120 Colliding balls:** PSSC Physics. **134 Ion drive:** NASA. **135 Halley's comet:** W. Liller. **135 Nucleus of Halley's comet:** European Space Agency. **138 Colliding galaxies:** NASA. **142 Wrench:** PSSC Physics. **142 Highjumper:** Courtesy of Dunia Young. **149 Air bag:** DaimlerChrysler AG, CC-BY-SA licensed.. **186 Nimitz freeway:** Courtesy of U.C. Berkeley Earth Sciences and Map Library. **191 Descartes:** French postal stamp. **205 Greyhound:** Line art by the author, based on a photo by Alex Lapuerta, CC-BY licensed. **208 Solar sail (artist's rendering):** Wikipedia user Paranoid, CC-BY-SA licensed. **218 Breaking trail:** Art by Walter E. Bohl. Image courtesy of the University of Michigan Museum of Art/School of Information and Library Studies. **224 Football player and old lady:** Hazel Abaya. **227 Biplane:** Open Clip Art Library, public domain. **236 Spider oscillations:** Emile, Le Floch, and Vollrath, *Nature* 440:621 (2006). **239 Runner:** Line art by B. Crowell, CC-BY-SA licensed. Based on a photo by Wikimedia Commons user Fengalon, public domain. **239 Rock climber:** Line art by B. Crowell, CC-BY-SA licensed. Based on a photo by Richard Peter/Deutsche Fotothek, CC-BY-SA licensed.. **241 ISS:** NASA/Crew of STS-132, public domain. **241 Skee ball:** Photo by Wikipedia user Joyous!, CC-BY-SA. **243 Rock climber:** Redrawn from a photo by Joseff Thomas, CC-BY. **251 High jump:** Thomas Eakins, public domain. **293 Explorer I:** NASA/JPL, public domain. **313 Hot air balloon:** Randy Oostdyk, CC-BY-SA licensed. **321 Carnot:** contemporary. **328 Space junk:** STK-generated images courtesy of CSSI (www.centerforspace.com). **329 Boltzmann's tomb:** Wikipedia user Daderot, CC-BY-SA licensed. **335 Difluoroethane molecule:** Wikipedia user Edgar181, public domain. **341 Otto cycle:** based on an animation by Wikipedia user UtzOnBike, CC-BY-SA licensed. **353 Electric bass:** Brynjar Vik, CC-BY license. **359 Superposition of pulses:** Photo from PSSC Physics. **370 Mount Wilson:** Andrew Dunn, cc-by-sa licensed. **372 X-15 shock wave:** NASA, public

domain. **389 Pan pipes:** Wikipedia user Andrew Dunn, CC-BY-SA licensed. **389 Flute:** Wikipedia user Grendelkhan, CC-BY-SA licensed. **390 Traffic:** Wikipedia user Diliff, CC-BY licensed. **397 GPS:** Wikipedia user HawaiianMama, CC-BY license. **397 Atomic clock on plane:** Copyright 1971, Associated press, used under U.S. fair use exception to copyright law. **398 Football pass:** Wikipedia user RMelon, CC-BY-SA licensed. **400 Horse:** From a public-domain photo by Eadweard Muybridge, 1872. **400 Satellite:** From a public-domain artist's conception of a GPS satellite, product of NASA. **401 Joan of Arc holding banner:** Ingres, 1854. **401 Joan of Arc interrogated:** Delaroche, 1856. **409 Muon storage ring at CERN:** (c) 1974 by CERN; used here under the U.S. fair use doctrine. **409 Colliding nuclei:** courtesy of RHIC. **416 Machine gunner's body:** Redrawn from a public-domain photo by Cpl. Sheila Brooks. **416 Machine gunner's head:** Redrawn from a sketch by Wenceslas Hollar, 17th century. **434 Eclipse:** 1919, public domain. **434 Newspaper headline:** 1919, public domain. **436 Photo of PET scanner:** Wikipedia user Hg6996, public domain. **436 Ring of detectors in PET scanner:** Wikipedia user Damato, public domain. **436 PET body scan:** Jens Langner, public domain. **473 Lightning:** C. Clark, NOAA photo library. **485 Millikan:** contemporary. **489 Thomson:** Harper's Monthly, 1904. **497 nuclear fuel pellets:** US DOE, public domain. **509 nuclear power plant:** Wikipedia user Stefan Kuhn, CC-BY-SA licensed. **515 GAMMASPHERE:** Courtesy of C.J. Lister and R.V.F. Janssens. **515 H bomb test:** public domain product of US DOE, Ivy Mike test. **515 Fatu Hiva Rainforest:** Wikipedia user Makemake, CC-BY-SA licensed. **515 fusion reactor:** "These images may be used free of charge for educational purposes but please use the acknowledgement 'photograph courtesy of EFDA-JET'". **515 sun:** SOHO (ESA and NASA). **518 Chernobyl map:** CIA Handbook of International Economic Statistics, 1996, public domain. **519 Fifty-foot woman:** Public domain due to nonrenewal of copyright. **519 Zombies:** Public domain due to an error by the distributor in failing to place a copyright notice on the film. **522 Horses:** (c) 2004 Elena Filatova. **522 Polar bear:** U.S. Fish and Wildlife Service, public domain. **523 UNILAC:** Copyrighted, not covered by the book's copyleft. **532 Knifefish:** Courtesy of Greg DeGreef. **543 Superconducting accelerator segment:** Courtesy of Argonne National Laboratory, managed and operated by the University of Chicago for the U.S. Department of Energy under contract no. W-31-109-ENG-38. **584 LIGO:** California Institute of Technology. **595 Topographical map:** United States Geological Survey, 19th century, uncopyrighted.. **614 Capacitors:** Wikipedia user de:Benutzer:Honina, CC-BY-SA licensed. **614 Inductors:** Wikipedia user de:Benutzer:Honina, CC-BY-SA licensed. **623 Ballasts:** Magnetic ballast: Wikimedia Commons use Atlant, CC-BY; solid-state ballast: Wikipedia user Anton, CC-BY-SA. **688 C.S. Wu:** Smithsonian Institution, believed to be public domain. **688 Swan Lake:** Peter Gerstbach, GFDL 1.2. **714 Faraday banknote:** fair use. **715 Ascending and Descending:** (c) 1960, M.C. Escher. **740 Laminated core:** Wikipedia user ArnoldReinhold, CC-BY-SA licensed. **740 Ferrite bead:** Photo of clip-on bead by the author; photo of built-in bead by Wikipedia user Stwalkerster, CC-BY-SA. **742 Levitating frog:** "Permission granted for this photo to be licensed under the GNU-type license by Lijnis Nelemans, High Field Magnet Laboratory, Radboud University Nijmegen.". **744 Core memory:** H.J. Sommer III, Professor of Mechanical Engineering, Penn State University, CC-BY. **744 Hysteresis curve:** Based on a figure by Wikipedia user Omegatron, CC-BY-SA licensed. **745 Fluxgate compass:** Wikipedia user Mike1024, public domain. **765 Rays of sunlight:** Wikipedia user PiccoloNamek, CC-BY-SA. **768 Jupiter and Io:** NASA/JPL/University of Arizona. **777 Ray-traced image:** Gilles Tran, Wikimedia Commons, public domain. **781 Praxinoscope:** Thomas B. Greenslade, Jr.. **783 The Sleeping Gypsy:** H. Rousseau, 1897. **786 Flower:** Based on a photo by Wikimedia Commons user Fir0002, CC-BY-SA. **786 Moon:** Wikimedia commons image. **787 Ladybug:** Redrawn from a photo by Wikimedia Commons user Gilles San Martin, CC-BY-SA. **799 Fish-eye lens:** Martin D"urrschnabel, CC-BY-SA. **801 Hubble space telescope:** NASA, public domain. **802 Cross-section of eye:** NEI. **802 Eye's anatomy:** After a public-domain drawing from NEI. **806 Water wave refracting:** Original photo from PSSC. **807 Ulcer:** Wikipedia user Aspersions, CC-BY-SA. **814 Diffraction of water waves:** Assembled from photos in PSSC. **814 Counterfactual lack of diffraction of water waves:** Assembled from photos in PSSC. **814 Diffraction of water waves:** Assembled from photos in PSSC. **815 Scaling of diffraction:** Assembled from photos in PSSC. **816 Huygens:** Contemporary painting?. **817 Diffraction of water waves:** Assembled from photos in PSSC. **818 Young:** Wikimedia Commons, "After a portrait by Sir Thomas Lawrence, From: Arthur Shuster & Arthur E. Shipley: Britain's Heritage of Science. London, 1917". **823 Single-slit diffraction of water waves:** PSSC. **823 Simulation of a single slit using three sources:** PSSC. **824 Pleiades:** NASA/ESA/AURA/Caltech, public domain. **824 Radio telescope:** Wikipedia user Hajor, CC-BY-SA. **826 Air wedge:** Franklin D. Jones, 1920, public domain. **841 Pleiades:** NASA/ESA/AURA/Caltech, public domain. **843 Anamorphic image:** Wikipedia user Istvan Orosz, CC-BY. **857 Mount St. Helens:** Public-domain image by Austin Post, USGS. **872 Ozone maps:** NASA/GSFC TOMS Team. **873 Photon interference photos:** Courtesy of Prof. Lyman Page. **914 Dance of the baby swans from Swan Lake:** Line art by B. Crowell, CC-BY-SA licensed. Based on a photo by Paata Vardanashvili, CC-BY licensed. **914 Circle dancing:** Franz von Stuck, 1910, public domain. **919 Superposition of pulses:** Photo from PSSC Physics. **941 Hindenburg:** Public domain product of the U.S. Navy. **960 Stern-Gerlach photo:** Gerlach's photo from 1922, public domain. **1000 Lissajous figure:** Wikimedia Commons user Alessio Damato,

CC-BY-SA.

Hints

Hints for chapter 2

Page 122, problem 16:

You can use either the chain-rule technique from page 83 or the technique prescribed in problem 15 on p. 122. The positions and velocities of the two masses are related to each other, and you'll need to use this relationship to eliminate one mass's position and velocity and get everything in terms of the other mass's position and velocity. The relationship between the two positions will involve some extraneous variables like the length of the string, which won't have any effect on your final result.

Page 122, problem 17:

This is similar to problem 16, but you're trying to find the combination of masses that will result in *zero* acceleration. In this problem, the distance dropped by one weight is different from, but still related to, the distance by which the other weight rises. Try relating the heights of the two weights to each other, so you can get the total gravitational energy in terms of only one of these heights.

Page 122, problem 18:

This is similar to problem 17, in that you're looking for a setup that will give zero acceleration, and the distance the middle weight rises or falls is not the same as the distance the other two weights fall or rise. The simplest approach is to get the three heights in terms of θ , so that you can write the gravitational energy in terms of θ .

Page 122, problem 19:

This is very similar to problems 16 and 17.

Page 122, problem 20:

The first two parts can be done more easily by setting $a = 1$, since the value of a only changes the distance scale. One way to do part b is by graphing.

Page 123, problem 22:

The condition for a circular orbit contains three unknowns, v , g , and r , so you can't just solve it for r . You'll need more equations to make three equations in three unknowns. You'll need a relationship between g and r , and also a relationship between v and r that uses the given fact that it's supposed to take 24 hours for an orbit.

Page 123, problem 25:

What does the total energy have to be if the projectile's velocity is exactly escape velocity? Write down conservation of energy, change v to dr/dt , separate the variables, and integrate.

Page 123, problem 26:

The analytic approach is a little cumbersome, although it can be done by using approximations like $1/\sqrt{1+\epsilon} \approx 1 - (1/2)\epsilon$. A more straightforward, brute-force method is simply to write a computer program that calculates U/m for a given point in spherical coordinates. By trial and error, you can fairly rapidly find the r that gives a desired value of U/m .

Page 125, problem 33:

Use calculus to find the minimum of U .

Page 125, problem 35:

The spring constant of this spring, k , is *not* the quantity you need in the equation for the period. What you need in that equation is the second derivative of the spring's energy with respect to

the position of the thing that's oscillating. You need to start by finding the energy stored in the spring as a function of the vertical position, y , of the mass. This is similar to example 23 on page 118.

Page 126, problem 37:

The variables x_1 and x_2 will adjust themselves to reach an equilibrium. Write down the total energy in terms of x_1 and x_2 , then eliminate one variable, and find the equilibrium value of the other. Finally, eliminate both x_1 and x_2 from the total energy, getting it just in terms of b .

Hints for chapter 3

Page 225, problem 20:

Write down two equations, one for Newton's second law applied to each object. Solve these for the two unknowns T and a .

Page 229, problem 41:

The whole expression for the amplitude has maxima where the stuff inside the square root is at a minimum, and vice versa, so you can save yourself a lot of work by just working on the stuff inside the square root. For normal, large values of Q , there are two extrema, one at $\omega = 0$ and one at resonance; one of these is a maximum and one is a minimum. You want to find out at what value of Q the zero-frequency extremum switches over from being a maximum to being a minimum.

Page 234, problem 69:

You can use the geometric interpretation of the dot product.

Page 235, problem 70:

The easiest way to do this problem is to use two different coordinate systems: one that's tilted to coincide with the upper slope, and one that's tilted to coincide with the lower one.

Hints for chapter 4

Page 294, problem 8:

The choice of axis theorem only applies to a closed system, or to a system acted on by a total force of zero. Even if the box is not going to rotate, its center of mass is going to accelerate, and this can still cause a change in its angular momentum, unless the right axis is chosen. For example, if the axis is chosen at the bottom right corner, then the box will start accumulating clockwise angular momentum, even if it is just accelerating to the right without rotating. Only by choosing the axis at the center of mass (or at some other point on the same horizontal line) do we get a constant, zero angular momentum.

Page 295, problem 11:

There are four forces on the wheel at first, but only three when it lifts off. Normal forces are always perpendicular to the surface of contact. Note that the corner of the step cannot be perfectly sharp, so the surface of contact for this force really coincides with the surface of the wheel.

Page 301, problem 35:

You'll need the result of problem 19 in order to relate the energy and angular momentum of a rigidly rotating body. Since this relationship involves a variable raised to a power, you can't just graph the data and get the moment of inertia directly. One way to get around this is to manipulate one of the variables to make the graph linear. Here is an example of this technique from another context. Suppose you were given a table of the masses, m , of cubical pieces of wood, whose sides had various lengths, b . You want to find a best-fit value for the density of

the wood. The relationship is $m = \rho b^3$. The graph of m versus b would be a curve, and you would not have any easy way to get the density from such a graph. But by graphing m versus b^3 , you can produce a graph that is linear, and whose slope equals the density.

Hints for chapter 6

Page 392, problem 4:

How could you change the values of x and t so that the value of y would remain the same? What would this represent physically?

Page 393, problem 8:

(a) The most straightforward approach is to apply the equation $\partial^2 y / \partial t^2 = (T/\mu) \partial^2 y / \partial x^2$. Although this equation was developed in the main text in the context of a straight string with a curvy wave on it, it works just as well for a circular loop; the left-hand side is simply the inward acceleration of any point on the rope. Note, however, that we've been assuming the string was (at least approximately) parallel to the x axis, which will only be true if you choose a specific value of x . You need to get an equation for y in terms of x in order to evaluate the right-hand side.

Page 394, problem 12:

The answers to the two parts are not the same.

Hints for chapter 7

Page 463, problem 28:

Apply the equivalence principle.

Hints for chapter 8

Page 527, problem 15:

The force on the lithium ion is the vector sum of all the forces of all the quadrillions of sodium and chlorine atoms, which would obviously be too laborious to calculate. Nearly all of these forces, however, are canceled by a force from an ion on the opposite side of the lithium.

Hints for chapter 9

Page 569, problem 20:

The approach is similar to the one used for the other problem, but you want to work with voltage and electrical energy rather than force.

Hints for chapter 10

Page 660, problem 15:

Use the approximation $(1 + \epsilon)^p \approx 1 + p\epsilon$, which is valid for small ϵ .

Page 663, problem 25:

First find the energy stored in a spherical shell extending from r to $r + dr$, then integrate to find the total energy.

Page 663, problem 26:

Since we have $t \ll r$, the volume of the membrane is essentially the same as if it was unrolled and flattened out, and the field's magnitude is nearly constant.

Page 664, problem 31:

The math is messy if you put the origin of your polar coordinates at the center of the disk. It comes out much simpler if you put the origin at the edge, right on top of the point at which we're trying to compute the voltage.

Page 665, problem 37:

There are various ways of doing this, but one easy and natural approach is to change the base of the exponent to e using the same method that we would use for real numbers.

Hints for chapter 11**Page 752, problem 24:**

A stable system has low energy; energy would have to be added to change its configuration.

Page 756, problem 41:

We're ignoring the fact that the light consists of little wavepackets, and imagining it as a simple sine wave. But wait, there's more good news! The energy density depends on the squares of the fields, which means the squares of some sine waves. Well, when you square a sine wave that varies from -1 to $+1$, you get a sine wave that goes from 0 to $+1$, and the average value of that sine wave is $1/2$. That means you don't have to do an integral like $U = \int (dU/dV) dV$. All you have to do is throw in the appropriate factor of $1/2$, and you can pretend that the fields have their constant values $\tilde{\mathbf{E}}$ and $\tilde{\mathbf{B}}$ everywhere.

Page 756, problem 42:

Use Faraday's law, and choose an Ampèrian surface that is a disk of radius R sandwiched between the plates.

Page 758, problem 51:

(a) Magnetic fields are created by currents, so once you've decided how currents behave under time-reversal, you can figure out how magnetic fields behave.

Hints for chapter 12**Page 844, problem 60:**

Expand $\sin \theta$ in a Taylor series around $\theta = 90^\circ$.

Solutions to selected problems

Solutions for chapter 0**Page 47, problem 6:**

$$134 \text{ mg} \times \frac{10^{-3} \text{ g}}{1 \text{ mg}} \times \frac{10^{-3} \text{ kg}}{1 \text{ g}} = 1.34 \times 10^{-4} \text{ kg}$$

Page 47, problem 8:

(a) Let's do 10.0 g and 1000 g . The arithmetic mean is 505 grams . It comes out to be 0.505 kg , which is consistent. (b) The geometric mean comes out to be 100 g or 0.1 kg , which is consistent. (c) If we multiply meters by meters, we get square meters. Multiplying grams by grams should give square grams! This sounds strange, but it makes sense. Taking the square root of square grams (g^2) gives grams again. (d) No. The superduper mean of two quantities with units of grams wouldn't even be something with units of grams! Related to this shortcoming is the fact that the superduper mean would fail the kind of consistency test carried out in the first two parts of the problem.

Page 48, problem 12:

(a) They're all defined in terms of the ratio of side of a triangle to another. For instance, the tangent is the length of the opposite side over the length of the adjacent side. Dividing meters

by meters gives a unitless result, so the tangent, as well as the other trig functions, is unitless.

(b) The tangent function gives a unitless result, so the units on the right-hand side had better cancel out. They do, because the top of the fraction has units of meters squared, and so does the bottom.

Page 49, problem 17:

$\Delta x = \frac{1}{2}at^2$, so for a fixed value of Δx , we have $t \propto 1/\sqrt{a}$. Translating this into the language of ratios gives $t_M/t_E = \sqrt{a_E/a_M} = \sqrt{3} = 1.7$.

Page 50, problem 19:

(a) Solving for $\Delta x = \frac{1}{2}at^2$ for a , we find $a = 2\Delta x/t^2 = 5.51 \text{ m/s}^2$. (b) $v = \sqrt{2a\Delta x} = 66.6 \text{ m/s}$.
 (c) The actual car's final velocity is less than that of the idealized constant-acceleration car. If the real car and the idealized car covered the quarter mile in the same time but the real car was moving more slowly at the end than the idealized one, the real car must have been going faster than the idealized car at the beginning of the race. The real car apparently has a greater acceleration at the beginning, and less acceleration at the end. This make sense, because every car has some maximum speed, which is the speed beyond which it cannot accelerate.

Page 51, problem 31:

$$1 \text{ mm}^2 \times \left(\frac{1 \text{ cm}}{10 \text{ mm}} \right)^2 = 10^{-2} \text{ cm}^2$$

Page 51, problem 32:

The bigger scope has a diameter that's ten times greater. Area scales as the square of the linear dimensions, so $A \propto d^2$, or in the language of ratios $A_1/A_2 = (d_1/d_2)^2 = 100$. Its light-gathering power is a hundred times greater.

Page 51, problem 33:

Since they differ by two steps on the Richter scale, the energy of the bigger quake is 10^4 times greater. The wave forms a hemisphere, and the surface area of the hemisphere over which the energy is spread is proportional to the square of its radius, $A \propto r^2$, or $r \propto \sqrt{A}$, which means $r_1/r_2 = \sqrt{A_1/A_2}$. If the amount of vibration was the same, then the surface areas must be in the ratio $A_1/A_2 = 10^4$, which means that the ratio of the radii is 10^2 .

Page 52, problem 38:

The cone of mixed gin and vermouth is the same shape as the cone of vermouth, but its linear dimensions are doubled. Translating the proportionality $V \propto L^3$ into an equation about ratios, we have $V_1/V_2 = (L_1/L_2)^3 = 8$. Since the ratio of the whole thing to the vermouth is 8, the ratio of gin to vermouth is 7.

Page 52, problem 40:

The proportionality $V \propto L^3$ can be restated in terms of ratios as $V_1/V_2 = (L_1/L_2)^3 = (1/10)^3 = 1/1000$, so scaling down the linear dimensions by a factor of 1/10 reduces the volume by 1/1000, to a milliliter.

Page 53, problem 41:

Let's estimate the Great Wall's volume, and then figure out how many bricks that would represent. The wall is famous because it covers pretty much all of China's northern border, so let's say it's 1000 km long. From pictures, it looks like it's about 10 m high and 10 m wide, so the total volume would be $10^6 \text{ m} \times 10 \text{ m} \times 10 \text{ m} = 10^8 \text{ m}^3$. If a single brick has a volume of 1 liter, or 10^{-3} m^3 , then this represents about 10^{11} bricks. If one person can lay 10 bricks in an hour

(taking into account all the preparation, etc.), then this would be 10^{10} man-hours.

Page 53, problem 44:

Directly guessing the number of jelly beans would be like guessing volume directly. That would be a mistake. Instead, we start by estimating the linear dimensions, in units of beans. The contents of the jar look like they're about 10 beans deep. Although the jar is a cylinder, its exact geometrical shape doesn't really matter for the purposes of our order-of-magnitude estimate. Let's pretend it's a rectangular jar. The horizontal dimensions are also something like 10 beans, so it looks like the jar has about $10 \times 10 \times 10$ or $\sim 10^3$ beans inside.

Solutions for chapter 1

Page 71, problem 12:

To the person riding the moving bike, bug A is simply going in circles. The only difference between the motions of the two wheels is that one is traveling through space, but motion is relative, so this doesn't have any effect on the bugs. It's equally hard for each of them.

Solutions for chapter 2

Page 120, problem 1:

(a) The energy stored in the gasoline is being changed into heat via frictional heating, and also probably into sound and into energy of water waves. Note that the kinetic energy of the propeller and the boat are not changing, so they are not involved in the energy transformation. (b) The cruising speed would be greater by a factor of the cube root of 2, or about a 26% increase.

Page 120, problem 2:

We don't have actual masses and velocities to plug in to the equation, but that's OK. We just have to reason in terms of ratios and proportionalities. Kinetic energy is proportional to mass and to the square of velocity, so B's kinetic energy equals $(13.4 \text{ J})(3.77)/(2.34)^2 = 9.23 \text{ J}$.

Page 120, problem 3:

Room temperature is about 20°C . The fraction of the energy that actually goes into heating the water is

$$\frac{(250 \text{ g})/(0.24 \text{ g}\cdot{}^\circ\text{C}/\text{J}) \times (100^\circ\text{C} - 20^\circ\text{C})}{(1.25 \times 10^3 \text{ J/s})(126 \text{ s})} = 0.53$$

So roughly half of the energy is wasted. The wasted energy might be in several forms: heating of the cup, heating of the oven itself, or leakage of microwaves from the oven.

Page 120, problem 5:

$$\begin{aligned} E_{total,i} &= E_{total,f} \\ PE_i + \text{heat}_i &= PE_f + KE_f + \text{heat}_f \\ \frac{1}{2}mv^2 &= PE_i - PE_f + \text{heat}_i - \text{heat}_f \\ &= -\Delta PE - \Delta \text{heat} \\ v &= \sqrt{2 \left(\frac{-\Delta PE - \Delta \text{heat}}{m} \right)} \\ &= 6.4 \text{ m/s} \end{aligned}$$

Solutions for chapter 3

Page 222, problem 4:

A conservation law is about addition: it says that when you add up a certain thing, the total always stays the same. Funkosity would violate the additive nature of conservation laws, because a two-kilogram mass would have twice as much funkosity as a pair of one-kilogram masses moving at the same speed.

Page 223, problem 12:

Momentum is a vector. The total momentum of the molecules is always zero, since the momenta in different directions cancel out on the average. Cooling changes individual molecular momenta, but not the total.

Page 224, problem 15:

$a = \Delta v / \Delta t$, and also $a = F/m$, so

$$\begin{aligned}\Delta t &= \frac{\Delta v}{a} \\ &= \frac{m\Delta v}{F} \\ &= \frac{(1000 \text{ kg})(50 \text{ m/s} - 20 \text{ m/s})}{3000 \text{ N}} \\ &= 10 \text{ s}\end{aligned}$$

Page 225, problem 23:

- (a) This is a measure of the box's resistance to a change in its state of motion, so it measures the box's mass. The experiment would come out the same in lunar gravity.
- (b) This is a measure of how much gravitational force it feels, so it's a measure of weight. In lunar gravity, the box would make a softer sound when it hit.
- (c) As in part a, this is a measure of its resistance to a change in its state of motion: its mass. Gravity isn't involved at all.

Page 228, problem 34:

- (a) The swimmer's acceleration is caused by the water's force on the swimmer, and the swimmer makes a backward force on the water, which accelerates the water backward.
- (b) The club's normal force on the ball accelerates the ball, and the ball makes a backward normal force on the club, which decelerates the club.
- (c) The bowstring's normal force accelerates the arrow, and the arrow also makes a backward normal force on the string. This force on the string causes the string to accelerate less rapidly than it would if the bow's force was the only one acting on it.
- (d) The tracks' backward frictional force slows the locomotive down. The locomotive's forward frictional force causes the whole planet earth to accelerate by a tiny amount, which is too small to measure because the earth's mass is so great.

Page 228, problem 37:

- (a) Spring constants in parallel add, so the spring constant has to be proportional to the cross-sectional area. Two springs in series give half the spring constant, three springs in series give $1/3$, and so on, so the spring constant has to be inversely proportional to the length. Summarizing, we have $k \propto A/L$.

- (b) With the Young's modulus, we have $k = (A/L)E$. The spring constant has units of N/m, so the units of E would have to be N/m².

Page 230, problem 44:

By conservation of momentum, the total momenta of the pieces after the explosion is the same

as the momentum of the firework before the explosion. However, there is no law of conservation of kinetic energy, only a law of conservation of energy. The chemical energy in the gunpowder is converted into heat and kinetic energy when it explodes. All we can say about the kinetic energy of the pieces is that their total is greater than the kinetic energy before the explosion.

Page 230, problem 45:

Let m be the mass of the little puck and $M = 2.3m$ be the mass of the big one. All we need to do is find the direction of the total momentum vector before the collision, because the total momentum vector is the same after the collision. Given the two components of the momentum vector $p_x = Mv$ and $p_y = mv$, the direction of the vector is $\tan^{-1}(p_y/p_x) = 23^\circ$ counterclockwise from the big puck's original direction of motion.

Page 233, problem 62:

We want to find out about the velocity vector \mathbf{v}_{BG} of the bullet relative to the ground, so we need to add Annie's velocity relative to the ground \mathbf{v}_{AG} to the bullet's velocity vector \mathbf{v}_{BA} relative to her. Letting the positive x axis be east and y north, we have

$$\begin{aligned} v_{BA,x} &= (140 \text{ mi/hr}) \cos 45^\circ \\ &= 100 \text{ mi/hr} \\ v_{BA,y} &= (140 \text{ mi/hr}) \sin 45^\circ \\ &= 100 \text{ mi/hr} \end{aligned}$$

and

$$\begin{aligned} v_{AG,x} &= 0 \\ v_{AG,y} &= 30 \text{ mi/hr.} \end{aligned}$$

The bullet's velocity relative to the ground therefore has components

$$v_{BG,x} = 100 \text{ mi/hr}$$

and

$$v_{BG,y} = 130 \text{ mi/hr.}$$

Its speed on impact with the animal is the magnitude of this vector

$$\begin{aligned} |\mathbf{v}_{BG}| &= \sqrt{(100 \text{ mi/hr})^2 + (130 \text{ mi/hr})^2} \\ &= 160 \text{ mi/hr} \end{aligned}$$

(rounded off to two significant figures).

Page 233, problem 63:

Since its velocity vector is constant, it has zero acceleration, and the sum of the force vectors acting on it must be zero. There are three forces acting on the plane: thrust, lift, and gravity. We are given the first two, and if we can find the third we can infer the plane's mass. The sum of the y components of the forces is zero, so

$$\begin{aligned} 0 &= F_{thrust,y} + F_{lift,y} + F_{g,y} \\ &= |\mathbf{F}_{thrust}| \sin \theta + |\mathbf{F}_{lift}| \cos \theta - mg. \end{aligned}$$

The mass is

$$\begin{aligned} m &= (|\mathbf{F}_{thrust}| \sin \theta + |\mathbf{F}_{lift}| \cos \theta) / g \\ &= 7.0 \times 10^4 \text{ kg.} \end{aligned}$$

Page 234, problem 64:

(a) Since the wagon has no acceleration, the total forces in both the x and y directions must be zero. There are three forces acting on the wagon: T , \mathbf{F}_g , and the normal force from the ground, \mathbf{F}_n . If we pick a coordinate system with x being horizontal and y vertical, then the angles of these forces measured counterclockwise from the x axis are $90^\circ - \phi$, 270° , and $90^\circ + \theta$, respectively. We have

$$\begin{aligned} F_{x,total} &= T \cos(90^\circ - \phi) + F_g \cos(270^\circ) + F_n \cos(90^\circ + \theta) \\ F_{y,total} &= T \sin(90^\circ - \phi) + F_g \sin(270^\circ) + F_n \sin(90^\circ + \theta), \end{aligned}$$

which simplifies to

$$\begin{aligned} 0 &= T \sin \phi - F_n \sin \theta \\ 0 &= T \cos \phi - F_g + F_n \cos \theta. \end{aligned}$$

The normal force is a quantity that we are not given and do not wish to find, so we should choose it to eliminate. Solving the first equation for $F_n = (\sin \phi / \sin \theta)T$, we eliminate F_n from the second equation,

$$0 = T \cos \phi - F_g + T \sin \phi \cos \theta / \sin \theta$$

and solve for T , finding

$$T = \frac{F_g}{\cos \phi + \sin \phi \cos \theta / \sin \theta}$$

Multiplying both the top and the bottom of the fraction by $\sin \theta$, and using the trig identity for $\sin(\theta + \phi)$ gives the desired result,

$$T = \frac{\sin \theta}{\sin(\theta + \phi)} F_g s$$

(b) The case of $\phi = 0$, i.e. pulling straight up on the wagon, results in $T = F_g$: we simply support the wagon and it glides up the slope like a chair-lift on a ski slope. In the case of $\phi = 180^\circ - \theta$, T becomes infinite. Physically this is because we are pulling directly into the ground, so no amount of force will suffice.

Page 234, problem 65:

(a) If there was no friction, the angle of repose would be zero, so the coefficient of static friction, μ_s , will definitely matter. We also make up symbols θ , m and g for the angle of the slope, the mass of the object, and the acceleration of gravity. The forces form a triangle just like the one in example 68 on page 207, but instead of a force applied by an external object, we have static friction, which is less than $\mu_s F_n$. As in that example, $F_s = mg \sin \theta$, and $F_s < \mu_s F_n$, so

$$mg \sin \theta < \mu_s F_n.$$

From the same triangle, we have $F_n = mg \cos \theta$, so

$$mg \sin \theta < \mu_s mg \cos \theta.$$

Rearranging,

$$\theta < \tan^{-1} \mu_s.$$

(b) Both m and g canceled out, so the angle of repose would be the same on an asteroid.

Page 242, problem 88:

(a) Based on units, we must have $g = kG\lambda/y$, where k is a unitless universal constant.

(b) For the actual calculation, we have

$$\begin{aligned} g &= \int dg_y \\ &= G \int \frac{dm}{r^2} \cos \theta, \end{aligned}$$

where θ is the angle between the perpendicular and the \mathbf{r} vector. Then $dm = \lambda dx$, $\cos \theta = y/r$, and $r = \sqrt{x^2 + y^2}$, so

$$\begin{aligned} g &= G \int_{-\infty}^{\infty} \frac{\lambda dx}{x^2 + y^2} \cdot \frac{b}{\sqrt{x^2 + y^2}} \\ &= G\lambda y \int_{-\infty}^{\infty} (x^2 + y^2)^{-3/2} dx. \end{aligned}$$

Even though this has limits of integration, this is an indefinite integral because it contains the variable y . It's nicer to clean this up by doing a change of variable to the unitless quantity $u = x/y$, giving

$$g = \frac{G\lambda}{y} \int_{-\infty}^{\infty} (u^2 + 1)^{-3/2} du.$$

The definite integral is the sort of thing that sane people these days will do using computer software. It equals 2. The result for the field is

$$g = \frac{2G\lambda}{y}.$$

Solutions for chapter 4

Page 294, problem 1:

The pliers are not moving, so their angular momentum remains constant at zero, and the total torque on them must be zero. Not only that, but each half of the pliers must have zero total torque on it. This tells us that the magnitude of the torque at one end must be the same as that at the other end. The distance from the axis to the nut is about 2.5 cm, and the distance from the axis to the centers of the palm and fingers are about 8 cm. The angles are close enough to 90° that we can pretend they're 90 degrees, considering the rough nature of the other assumptions and measurements. The result is $(300 \text{ N})(2.5 \text{ cm}) = (F)(8 \text{ cm})$, or $F = 90 \text{ N}$.

Page 301, problem 37:

The foot of the rod is moving in a circle relative to the center of the rod, with speed $v = \pi b/T$, and acceleration $v^2/(b/2) = (\pi^2/8)g$. This acceleration is initially upward, and is greater in magnitude than g , so the foot of the rod will lift off without dragging. We could also worry about whether the foot of the rod would make contact with the floor again before the rod finishes up flat on its back. This is a question that can be settled by graphing, or simply by inspection of figure i on page 282. The key here is that the two parts of the acceleration are

both independent of m and b , so the result is universal, and it does suffice to check a graph in a single example. In practical terms, this tells us something about how difficult the trick is to do. Because $\pi^2/8 = 1.23$ isn't much greater than unity, a hit that is just a little too weak (by a factor of $1.23^{1/2} = 1.11$) will cause a fairly obvious qualitative change in the results. This is easily observed if you try it a few times with a pencil.

Page 301, problem 37:

The foot of the rod is moving in a circle relative to the center of the rod, with speed $v = \pi b/T$, and acceleration $v^2/(b/2) = (\pi^2/8)g$. This acceleration is initially upward, and is greater in magnitude than g , so the foot of the rod will lift off without dragging. We could also worry about whether the foot of the rod would make contact with the floor again before the rod finishes up flat on its back. This is a question that can be settled by graphing, or simply by inspection of figure i on page 282. The key here is that the two parts of the acceleration are both independent of m and b , so the result is universal, and it does suffice to check a graph in a single example. In practical terms, this tells us something about how difficult the trick is to do. Because $\pi^2/8 = 1.23$ isn't much greater than unity, a hit that is just a little too weak (by a factor of $1.23^{1/2} = 1.11$) will cause a fairly obvious qualitative change in the results. This is easily observed if you try it a few times with a pencil.

Page 303, problem 45:

The moment of inertia is $I = \int r^2 dm$. Let the ring have total mass M and radius b . The proportionality

$$\frac{M}{2\pi} = \frac{dm}{d\theta}$$

gives a change of variable that results in

$$I = \frac{M}{2\pi} \int_0^{2\pi} r^2 d\theta.$$

If we measure θ from the axis of rotation, then $r = b \sin \theta$, so this becomes

$$I = \frac{Mb^2}{2\pi} \int_0^{2\pi} \sin^2 \theta d\theta.$$

The integrand averages to $1/2$ over the 2π range of integration, so the integral equals π . We therefore have $I = \frac{1}{2}Mb^2$. This is, as claimed, half the value for rotation about the symmetry axis.

Solutions for chapter 5

Page 349, problem 11:

(a) We have

$$\begin{aligned} dP &= \rho g dy \\ \Delta P &= \int \rho g dy, \end{aligned}$$

and since we're taking water to be incompressible, and g doesn't change very much over 11 km

of height, we can treat ρ and g as constants and take them outside the integral.

$$\begin{aligned}\Delta P &= \rho g \Delta y \\ &= (1.0 \text{ g/cm}^3)(9.8 \text{ m/s}^2)(11.0 \text{ km}) \\ &= (1.0 \times 10^3 \text{ kg/m}^3)(9.8 \text{ m/s}^2)(1.10 \times 10^4 \text{ m}) \\ &= 1.0 \times 10^8 \text{ Pa} \\ &= 1.0 \times 10^3 \text{ atm.}\end{aligned}$$

The precision of the result is limited to a few percent, due to the compressibility of the water, so we have at most two significant figures. If the change in pressure were exactly a thousand atmospheres, then the pressure at the bottom would be 1001 atmospheres; however, this distinction is not relevant at the level of approximation we're attempting here.

(b) Since the air in the bubble is in thermal contact with the water, it's reasonable to assume that it keeps the same temperature the whole time. The ideal gas law is $PV = nkT$, and rewriting this as a proportionality gives

$$V \propto P^{-1},$$

or

$$\frac{V_f}{V_i} = \left(\frac{P_f}{P_i} \right)^{-1} \approx 10^3.$$

Since the volume is proportional to the cube of the linear dimensions, the growth in radius is about a factor of 10.

Page 349, problem 12:

(a) Roughly speaking, the thermal energy is $\sim k_B T$ (where k_B is the Boltzmann constant), and we need this to be on the same order of magnitude as ke^2/r (where k is the Coulomb constant). For this type of rough estimate it's not especially crucial to get all the factors of two right, but let's do so anyway. Each proton's average kinetic energy due to motion along a particular axis is $(1/2)k_B T$. If two protons are colliding along a certain line in the center-of-mass frame, then their average combined kinetic energy due to motion along that axis is $2(1/2)k_B T = k_B T$. So in fact the factors of 2 cancel. We have $T = ke^2/k_B r$.

(b) The units are $K = (J \cdot m/C^2)(C^2)/((J/K) \cdot m)$, which does work out.

(c) The numerical result is $\sim 10^{10}$ K, which as suggested is much higher than the temperature at the core of the sun.

Page 351, problem 13:

If the full-sized brick A undergoes some process, such as heating it with a blowtorch, then we want to be able to apply the equation $\Delta S = Q/T$ to either the whole brick or half of it, which would be identical to B. When we redefine the boundary of the system to contain only half of the brick, the quantities ΔS and Q are each half as big, because entropy and energy are additive quantities. T , meanwhile, stays the same, because temperature isn't additive — two cups of coffee aren't twice as hot as one. These changes to the variables leave the equation consistent, since each side has been divided by 2.

Page 351, problem 14:

(a) If the expression $1 + by$ is to make sense, then by has to be unitless, so b has units of m^{-1} . The input to the exponential function also has to be unitless, so k also has units of m^{-1} . The only factor with units on the right-hand side is P_o , so P_o must have units of pressure, or Pa.

(b)

$$\begin{aligned} dP &= \rho g dy \\ \rho &= \frac{1}{g} \frac{dP}{dy} \\ &= \frac{P_o}{g} e^{-ky} (-k - kby + b) \end{aligned}$$

(c) The three terms inside the parentheses on the right all have units of m^{-1} , so it makes sense to add them, and the factor in parentheses has those units. The units of the result from b then look like

$$\begin{aligned} \frac{\text{kg}}{\text{m}^3} &= \frac{\text{Pa}}{\text{m/s}^2} \text{m}^{-1} \\ &= \frac{\text{N/m}^2}{\text{m}^2/\text{s}^2} \\ &= \frac{\text{kg} \cdot \text{m}^{-1} \cdot \text{s}^{-2}}{\text{m}^2/\text{s}^2}, \end{aligned}$$

which checks out.

Solutions for chapter 7

Page 461, problem 17:

(a) Plugging in, we find

$$\sqrt{\frac{1-w}{1+w}} = \sqrt{\frac{1-u}{1+u}} \sqrt{\frac{1-v}{1+v}}.$$

(b) First let's simplify by squaring both sides.

$$\frac{1-w}{1+w} = \frac{1-u}{1+u} \cdot \frac{1-v}{1+v}.$$

For convenience, let's write A for the right-hand side of this equation. We then have

$$\begin{aligned} \frac{1-w}{1+w} &= A \\ 1-w &= A+Aw. \end{aligned}$$

Solving for w ,

$$\begin{aligned} w &= \frac{1-A}{1+A} \\ &= \frac{(1+u)(1+v)-(1-u)(1-v)}{(1+u)(1+v)+(1-u)(1-v)} \\ &= \frac{2(u+v)}{2(1+uv)} \\ &= \frac{u+v}{1+uv} \end{aligned}$$

(c) This is all in units where $c = 1$. The correspondence principle says that we should get $w \approx u + v$ when both u and v are small compared to 1. Under those circumstances, uv is the

product of two very small numbers, which makes it very, very small. Neglecting this term in the denominator, we recover the nonrelativistic result.

Page 461, problem 18:

Among the spacelike vectors, \mathbf{a} and \mathbf{e} are clearly congruent, because they're the same except for a rotation in space; this is the same as the definition of congruence in ordinary Euclidean geometry, where rotation doesn't matter. Vector \mathbf{b} is also congruent to these, since it represents an interval $3^2 - 5^2 = -4^2$, just like the other two.

The lightlike vectors \mathbf{c} and \mathbf{d} both represent intervals of zero, so they're congruent, even though \mathbf{c} is a double-scale version of \mathbf{d} .

The timelike vectors \mathbf{f} and \mathbf{g} are not congruent to each other or to any of the others; \mathbf{f} represents an interval of 2^2 , while \mathbf{g} 's interval is 4^2 .

Page 462, problem 22:

At the center of each of the large triangle's sides, the angles add up to 180° because they form a straight line. Therefore $4s = S + 3 \times 180^\circ$, so $S - 180^\circ = 4(s - 180^\circ)$.

Page 463, problem 28:

By the equivalence principle, we can adopt a frame tied to the tossed clock, B, and in this frame there is no gravitational field. We see a desk and clock A go by. The desk applies a force to clock A, decelerating it and then reaccelerating it so that it comes back. We've already established that the effect of motion is to slow down time, so clock A reads a smaller time interval.

Page 464, problem 32:

To make the units make sense, we need to make sure that both sides of the \approx sign have the same units, and also that both terms on the right-hand side have the same units. Everything is unitless except for the second term on the right, so we add a factor of c^{-2} to fix it:

$$\gamma \approx 1 + \frac{v^2}{2c^2}.$$

Solutions for chapter 9

Page 566, problem 1:

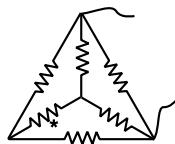
$$\Delta t = \Delta q/I = e/I = 0.16 \mu s$$

Page 567, problem 12:

In series, they give $11 \text{ k}\Omega$. In parallel, they give $(1/1 \text{ k}\Omega + 1/10 \text{ k}\Omega)^{-1} = 0.9 \text{ k}\Omega$.

Page 570, problem 25:

The actual shape is irrelevant; all we care about is what's connected to what. Therefore, we can draw the circuit flattened into a plane. Every vertex of the tetrahedron is adjacent to every other vertex, so any two vertices to which we connect will give the same resistance. Picking two arbitrarily, we have this:



This is unfortunately a circuit that cannot be converted into parallel and series parts, and that's what makes this a hard problem! However, we can recognize that by symmetry, there is zero current in the resistor marked with an asterisk. Eliminating this one, we recognize the

whole arrangement as a triple parallel circuit consisting of resistances R , $2R$, and $2R$. The resulting resistance is $R/2$.

Page 571, problem 29:

(a) Conservation of energy gives

$$\begin{aligned} U_A &= U_B + K_B \\ K_B &= U_A - U_B \\ \frac{1}{2}mv^2 &= e\Delta V \\ v &= \sqrt{\frac{2e\Delta V}{m}} \end{aligned}$$

(b) Plugging in numbers, we get 5.9×10^7 m/s. This is about 20% of the speed of light, so the nonrelativistic assumption was good to at least a rough approximation.

Page 572, problem 32:

It's much more practical to measure voltage differences. To measure a current, you have to break the circuit somewhere and insert the meter there, but it's not possible to disconnect the circuits sealed inside the board.

Solutions for chapter 10

Page 660, problem 16:

By symmetry, the field is always directly toward or away from the center. We can therefore calculate it along the x axis, where $r = x$, and the result will be valid for any location at that distance from the center. The electric field is minus the derivative of the potential,

$$\begin{aligned} E &= -\frac{dV}{dx} \\ &= -\frac{d}{dx}(x^{-1}e^{-x}) \\ &= x^{-2}e^{-x} + x^{-1}e^{-x} \end{aligned}$$

At small x , near the proton, the first term dominates, and the exponential is essentially 1, so we have $E \propto x^{-2}$, as we expect from the Coulomb force law. At large x , the second term dominates, and the field approaches zero faster than an exponential.

Page 668, problem 56:

$$\begin{aligned} \sin(a+b) &= (e^{i(a+b)} - e^{-i(a+b)})/2i \\ &= (e^{ia}e^{ib} - e^{-ia}e^{-ib})/2i \\ &= [(\cos a + i \sin a)(\cos b + i \sin b) - (\cos a - i \sin a)(\cos b - i \sin b)]/2i \\ &= \cos a \sin b + \sin a \cos b \end{aligned}$$

By a similar computation, we find $\cos(a+b) = \cos a \cos b - \sin a \sin b$.

Page 668, problem 57:

If $z^3 = 1$, then we know that $|z| = 1$, since cubing z cubes its magnitude. Cubing z triples its

argument, so the argument of z must be a number that, when tripled, is equivalent to an angle of zero. There are three possibilities: $0 \times 3 = 0$, $(2\pi/3) \times 3 = 2\pi$, and $(4\pi/3) \times 3 = 4\pi$. (Other possibilities, such as $(32\pi/3)$, are equivalent to one of these.) The solutions are:

$$z = 1, e^{2\pi i/3}, e^{4\pi i/3}$$

Page 668, problem 59:

We have $D = q\ell$ and $F_x = qb\ell = Db$. Since b is the same thing, in this example, as $\partial E_x / \partial x$, our equation for F is consistent with the result of example 7 on p. 591. and also, as claimed, depends on q and ℓ only via D .

Solutions for chapter 11

Page 758, problem 51:

- (a) For a material object, $\mathbf{p} = m\mathbf{v}$. The velocity vector reverses itself, but mass is still positive, so the momentum vector is reversed.
- (b) In the forward-time universe, conservation of momentum is $\mathbf{p}_{1,i} + \mathbf{p}_{2,i} = \mathbf{p}_{1,f} + \mathbf{p}_{2,f}$. In the backward-time universe, all the momenta are reversed, but that just negates both sides of the equation, so momentum is still conserved.

Page 758, problem 51:

- (a) For a material object, $\mathbf{p} = m\mathbf{v}$. The velocity vector reverses itself, but mass is still positive, so the momentum vector is reversed.
- (b) In the forward-time universe, conservation of momentum is $\mathbf{p}_{1,i} + \mathbf{p}_{2,i} = \mathbf{p}_{1,f} + \mathbf{p}_{2,f}$. In the backward-time universe, all the momenta are reversed, but that just negates both sides of the equation, so momentum is still conserved.

Page 760, problem 54:

Note that in the Biot-Savart law, the variable \mathbf{r} is defined as a vector that points from the current to the point at which the field is being calculated, whereas in the polar coordinates used to express the equation of the spiral, the vector more naturally points the opposite way. This requires some fiddling with signs, which I'll suppress, and simply identify $d\ell$ with $d\mathbf{r}$.

$$\mathbf{B} = \frac{kI}{c^2} \int \frac{d\ell \times \mathbf{r}}{r^3}$$

The vector $d\mathbf{r}$ has components $dx = w(\cos \theta - \theta \sin \theta)$ and $dy = w(\sin \theta + \theta \cos \theta)$. Evaluating the vector cross product, and substituting θ/w for r , we find

$$\begin{aligned} \mathbf{B} &= \frac{kI}{c^2 w} \int \frac{\theta(\cos \theta \sin \theta - \theta \sin^2 \theta - \cos \theta \sin \theta - \theta \cos^2 \theta) d\theta}{\theta^3} \\ &= \frac{kI}{c^2 w} \int \frac{d\theta}{\theta} \\ &= \frac{kI}{c^2 w} \ln \frac{\theta_2}{\theta_1} \\ &= \frac{kI}{c^2 w} \ln \frac{b}{a} \end{aligned}$$

Solutions for chapter 12

Page 829, problem 4:

Because the surfaces are flat, you get specular reflection. In specular reflection, all the reflected

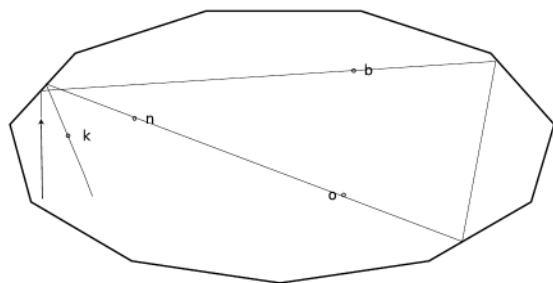
rays go in one direction. Unless the plane is directly overhead, that direction won't be the right direction to make the rays come back to the radar station.



This is different from a normal plane, which has complicated, bumpy surfaces. These surfaces give diffuse reflection, which spreads the reflected rays randomly in more or less every possible direction.

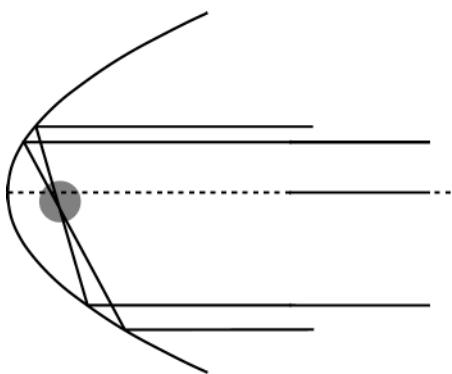
Page 829, problem 5:

It spells "bonk."



Page 830, problem 6:

- (a) The rays all cross at pretty much the same place, given the accuracy with which we can draw them.
- (b) It could be used to cook food, for instance. All the sunlight is concentrated in a small area.
- (c) Put the lightbulb at the point where the rays cross. The outgoing rays will then form a parallel beam going out to the right.



Page 831, problem 11:

The magnification is the ratio of the image's size to the object's size. It has nothing to do with the person's location. The angular magnification, however, does depend on the person's location, because things farther away subtend smaller angles. The distance to the actual object is not changed significantly, since it's zillions of miles away in outer space, but the distance to the image does change if the observer's point of view changes. If you can get closer to the image, the angular magnification is greater.

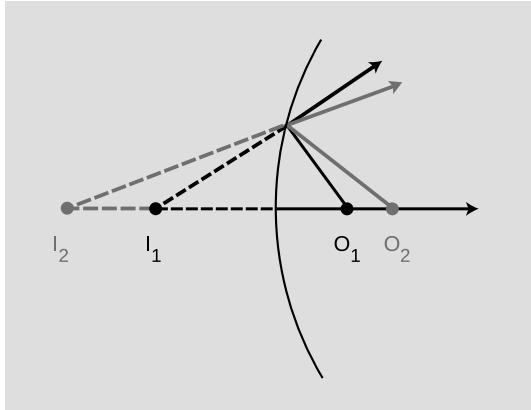
Page 832, problem 15:

For a flat mirror, d_i and d_o are equal, so the magnification is 1, i.e., the image is the same size

as the object.

Page 832, problem 16:

See the ray diagram below. Decreasing θ_o decreases θ_i , so the equation $\theta_f = \pm\theta_i + \pm\theta_o$ must have opposite signs on the right. Since θ_o is bigger than θ_i , the only way to get a positive θ_f is if the signs are $\theta_f = -\theta_i + \theta_o$. This gives $1/f = -1/d_i + 1/d_o$.



Page 832, problem 19:

- (a) The object distance is less than the focal length, so the image is virtual: because the object is so close, the cone of rays is diverging too strongly for the mirror to bring it back to a focus.
- (b) Now the object distance is greater than the focal length, so the image is real. (c),(d) A diverging mirror can only make virtual images.

Page 832, problem 20:

- (a) In problem #2 we found that the equation relating the object and image distances was of the form $1/f = -1/d_i + 1/d_o$. Let's make $f = 1.00$ m. To get a virtual image we need $d_o < f$, so let $d_o = 0.50$ m. Solving for d_i , we find $d_i = 1/(1/d_o - 1/f) = 1.00$ m. The magnification is $M = d_i/d_o = 2.00$. If we change d_o to 0.55 m, the magnification becomes 2.22. The magnification changes somewhat with distance, so the store's ad must be assuming you'll use the mirror at a certain distance. It can't have a magnification of 5 at all distances.
- (b) Theoretically yes, but in practical terms no. If you go through a calculation similar to the one in part a, you'll find that the images of both planets are formed at almost exactly the same d_i , $d_i = f$, since $1/d_o$ is pretty close to zero for any astronomical object. The more distant planet has an image half as big ($M = d_i/d_o$, and d_o is doubled), but we're talking about *angular* magnification here, so what we care about is the angular size of the image compared to the angular size of the object. The more distant planet has half the angular size, but its image has half the angular size as well, so the angular magnification is the same. If you think about it, it wouldn't make much sense for the angular magnification to depend on the planet's distance — if it did, then determining astronomical distances would be much easier than it actually is!

Page 832, problem 21:

- (a) This occurs when the d_i is infinite. Let's say it's a converging mirror creating a virtual image, as in problems 2 and 3. Then we'd get an infinite d_i if we put $d_o = f$, i.e., the object is at the focal point of the mirror. The image is infinitely large, but it's also infinitely far away, so its angular size isn't infinite; an angular size can never be more than about 180° since you can't see in back of your head!
- (b) It's not possible to make the magnification infinite by having $d_o = 0$. The image location and object location are related by $1/f = 1/d_o - 1/d_i$, so $1/d_i = 1/d_o - 1/f$. If d_o is zero, then

$1/d_o$ is infinite, $1/d_i$ is infinite, and d_i is zero as well. In other words, as d_o approaches zero, so does d_i , and d_i/d_o doesn't blow up. Physically, the mirror's curvature becomes irrelevant from the point of view of a tiny flea sitting on its surface: the mirror seems flat to the flea. So physically the magnification would be 1, not infinity, for very small values of d_o .

Page 834, problem 27:

The refracted ray that was bent closer to the normal in the plastic when the plastic was in air will be bent farther from the normal in the plastic when the plastic is in water. It will become a diverging lens.

Page 834, problem 29:

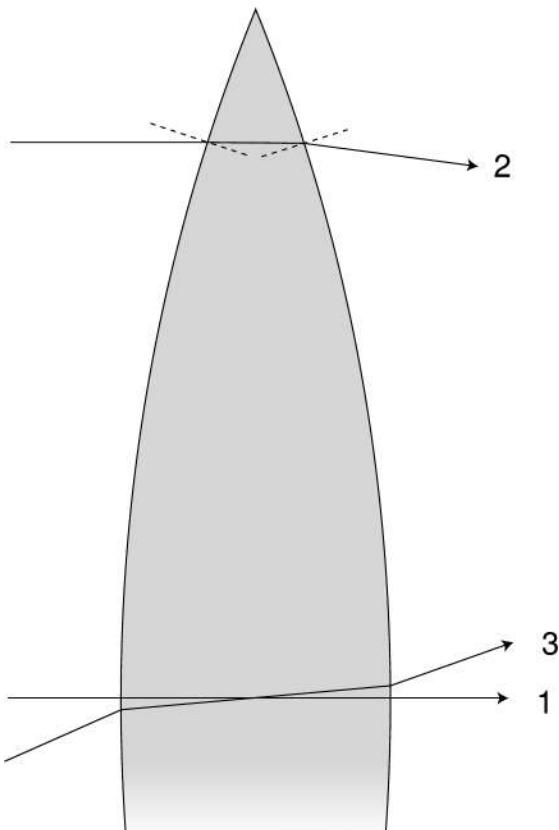
Refraction occurs only at the boundary between two substances, which in this case means the surface of the lens. Light doesn't get bent at all inside the lens, so the thickness of the lens isn't really what's important. What matters is the angles of the lens' surfaces at various points.

Ray 1 makes an angle of zero with respect to the normal as it enters the lens, so it doesn't get bent at all, and likewise at the back.

At the edge of the lens, 2, the front and back are not parallel, so a ray that traverses the lens at the edge ends up being bent quite a bit.

Although I drew both ray 1 and ray 2 coming in along the axis of the lens, it really doesn't matter. For instance, ray 3 bends on the way in, but bends an equal amount on the way out, so it still emerges from the lens moving in the same direction as the direction it originally had.

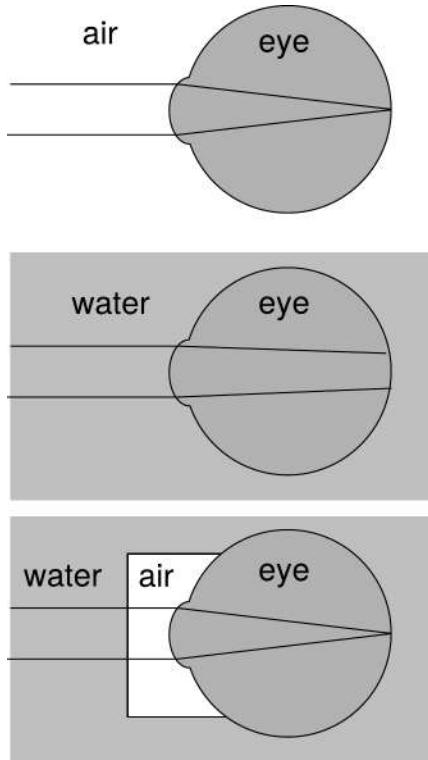
Summarizing and systematizing these observations, we can say that for a ray that enters the lens at the center, where the surfaces are parallel, the sum of the two deflection angles is zero. Since the total deflection is zero at the center, it must be larger away from the center.



Page 835, problem 31:

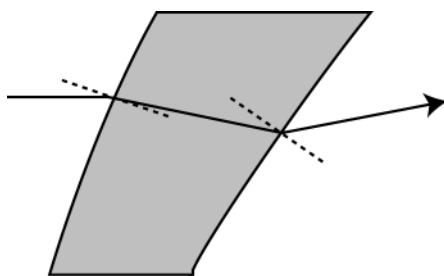
Normally, in air, your eyes do most of their focusing at the air-eye boundary. When you swim without goggles, there is almost no difference in speed at the water-eye interface, so light is not strongly refracted there (see figure), and the image is far behind the retina.

Goggles fix this problem for the following reason. The light rays cross a water-air boundary as they enter the goggles, but they're coming in along the normal, so they don't get bent. At the air-eye boundary, they get bent the same amount they normally would when you weren't swimming.



Page 835, problem 33:

(a) See the figure below. The first refraction clearly bends it inward. However, the back surface of the lens is more slanted, so the ray makes a bigger angle with respect to the normal at the back surface. The bending at the back surface is therefore greater than the bending at the front surface, and the ray ends up being bent *outward* more than inward.



(b) Lens 2 must act the same as lens 1. It's diverging. One way of knowing this is time-reversal symmetry: if we flip the original figure over and then reverse the direction of the ray, it's still a valid diagram.

Lens 3 is diverging like lens 1 on top, and diverging like lens 2 on the bottom. It's a diverging lens.

As for lens 4, any close-up diagram we draw of a particular ray passing through it will look exactly like the corresponding close-up diagram for some part of lens 1. Lens 4 behaves the same as lens 1.

Page 836, problem 39:

Since d_o is much greater than d_i , the lens-film distance d_i is essentially the same as f . (a) Splitting the triangle inside the camera into two right triangles, straightforward trigonometry gives

$$\theta = 2 \tan^{-1} \frac{w}{2f}$$

for the field of view. This comes out to be 39° and 64° for the two lenses. (b) For small angles, the tangent is approximately the same as the angle itself, provided we measure everything in radians. The equation above then simplifies to

$$\theta = \frac{w}{f}$$

The results for the two lenses are $.70 \text{ rad} = 40^\circ$, and $1.25 \text{ rad} = 72^\circ$. This is a decent approximation.

(c) With the 28-mm lens, which is closer to the film, the entire field of view we had with the 50-mm lens is now confined to a small part of the film. Using our small-angle approximation $\theta = w/f$, the amount of light contained within the same angular width θ is now striking a piece of the film whose linear dimensions are smaller by the ratio 28/50. Area depends on the square of the linear dimensions, so all other things being equal, the film would now be overexposed by a factor of $(50/28)^2 = 3.2$. To compensate, we need to shorten the exposure by a factor of 3.2.

Page 840, problem 48:

You don't want the wave properties of light to cause all kinds of funny-looking diffraction effects. You want to see the thing you're looking at in the same way you'd see a big object. Diffraction effects are most pronounced when the wavelength of the light is relatively large compared to the size of the object the light is interacting with, so red would be the worst. Blue light is near the short-wavelength end of the visible spectrum, which would be the best.

Page 840, problem 49:

- (a) You can tell it's a single slit because of the double-width central fringe.
- (b) Four fringes on the top pattern are about 23.5 mm, while five fringes on the bottom one are about 14.5 mm. The spacings are 5.88 and 2.90 mm, with a ratio of 2.03. A smaller d leads to larger diffraction angles, so the width of the slit used to make the bottom pattern was almost exactly twice as wide as the one used to make the top one.

Page 841, problem 51:

For the size of the diffraction blob, we have:

$$\begin{aligned} \frac{\lambda}{d} &\sim \sin \theta \\ &\approx \theta \\ \theta &\sim \frac{700 \text{ nm}}{10 \text{ m}} \\ &\approx 10^{-7} \text{ radians} \end{aligned}$$

For the actual angular size of the star, the small-angle approximation gives

$$\begin{aligned}\theta &\sim \frac{10^9 \text{ m}}{10^{17} \text{ m}} \\ &= 10^{-8} \text{ radians}\end{aligned}$$

The diffraction blob is ten times bigger than the actual disk of the star, so we can never make an image of the star itself in this way.

Page 841, problem 52:

(a) The patterns have two structures, a coarse one and a fine one. You can look up in the book which corresponds to w and which to d , or just use the fact that smaller features make bigger diffraction angles. The top and middle patterns have the same coarse spacing, so they have the same w . The fine structure in the top pattern has 7 fringes in 12.5 mm, for a spacing of 1.79 mm, while the middle pattern has 11 fringes in 41.5 mm, giving a spacing of 3.77 mm. The value of d for the middle pattern is therefore $(0.50 \text{ mm})(1.79/3.77) = 0.23 \text{ mm}$.

(b) This one has about the same d as the top one (it's difficult to measure accurately because each group has only a small number of fringes), but the coarse spacing is different, indicating a different value of w . It has two coarse groupings in 23 mm, i.e., a spacing of 12.5 mm. The coarse groupings in the original pattern were about 23 mm apart, so there is a factor of two between the $w = 0.04 \text{ mm}$ of the top pattern and the $w = 0.08 \text{ mm}$ of the bottom one.

Page 842, problem 55:

The equation, solved for θ , is $\theta = \sin^{-1}(m\lambda/d)$. The sine function only ranges from -1 to $+1$, so the inverse sine is undefined for $|m\lambda/d| > 1$, i.e., $|m| > d/\lambda$. Physically, we only get fringes out to angles of 90 degrees (the inverse sine of 1) on both sides, corresponding to values of m less than d/λ .

Page 844, problem 59:

One surface is curved outward and one inward. Therefore the minus sign applies in the lens-maker's equation. Since the radii of curvature are equal, the quantity $1/r_1 - 1/r_2$ equals zero, and the resulting focal length is infinite. A big focal length indicates a weak lens. An infinite focal length tells us that the lens is infinitely weak — it doesn't focus or defocus rays at all.

Page 845, problem 63:

We have $n = \sin \phi / \sin \theta$. Doing implicit differentiation, we find $dn = -\sin \phi (\cos \theta / \sin^2 \theta) d\theta$, which can be rewritten as $dn = -n \cot \theta d\theta$. This can be minimized by making θ as big as possible. To make θ as big as possible, we want ϕ to be as close as possible to 90 degrees, i.e., almost grazing the surface of the tank.

This result makes sense, because we're depending on refraction in order to get a measurement of n . At $\phi = 0$, we get $\theta = 0$, which provides no information at all about the index of refraction — the error bars become infinite. The amount of refraction increases as the angles get bigger.

Page 846, problem 64:

(a) The situation being described requires a real image, since the rays need to converge at a point on Becky's neck. See the ray diagram drawn with thick lines, showing object location o and image location i .



If we move the object farther away, to o' the cone of rays intercepted by the lens (thin lines) is less strongly diverging, and the lens is able to bring it to a closer focus, at i' . In the diagrams, we see that a smaller θ_o leads to a larger θ_i , so the signs in the equation $\pm\theta_o \pm \theta_i = \theta_f$ must be the same, and therefore both positive, since θ_f is positive by definition. The equation relating the image and object locations must be $1/f = 1/d_o + 1/d_i$.

(b) The case with $d_i = f$ is not possible, because then we need $1/d_o = 0$, i.e., $d_o = \infty$. Although it is possible in principle to have an object so far away that it is practically at infinity, that is not possible in this situation, since Zahra can't take her lens very far away from the fire. By the way, this means that the *focal length* f is not where the *focus* happens — the focus happens at d_i .

For similar reasons, we can't have $d_o = f$.

Since all the variables are positive, we must have $1/d_o$ and $1/d_i$ both less than $1/f$. This implies that $d_o > f$ and $d_i > f$. Of the nine logical possibilities in the table, only this one is actually possible for this real image.

Solutions for chapter 13

Page 952, problem 48:

The expressions $|\Psi|^2$ and $|\Psi^2|$ are identical, because the magnitude of a product is the product of the magnitudes. These expressions give positive real numbers as their results, which makes sense for a probability density. The expression Ψ^2 need not be real, and if it is real, it may be negative. It cannot be interpreted as a probability density. As a concrete example, suppose that $\Psi = bi$, where b is a real number with units. Then $|\Psi|^2 = |\Psi^2| = b^2$, which is real and positive, but $\Psi^2 = -b^2$, which clearly can't be interpreted probabilistically, because it's negative.

Page 952, problem 49:

(a) The quantity $x - y$ vanishes along the line $y = x$ lying in the first quadrant at a 45-degree angle between the axes. Squaring produces a trough parallel to this line, with a parabolic cross-section. Geometrically, the Laplacian can be interpreted as a measure of how much the value of f at a point differs from its average value on a small circle centered on that point. The trough is concave up, so we can predict that the Laplacian will be positive everywhere.

(b) The zero result is clearly wrong because it disagrees with our conclusion from part a that the Laplacian is positive. A correct calculation gives $\partial^2(x - y)^2 / \partial x^2 + \partial^2(x - y)^2 / \partial y^2 = 4$.

(c) If we rotate our coordinate axes counterclockwise by 45 degrees, then we have a parabolic trough oriented along the x axis. In terms of these new coordinates, $\partial f / \partial x = 0$, while $\partial f / \partial y$ is nonzero almost everywhere.

Remark: The mistake described in the question is a common one, and is apparently based on the idea that the notation ∇^2 must mean applying an operator ∇ twice. For those with some exposure to vector calculus, it may be of interest to note that the Laplacian *is* equivalent to the divergence of the gradient, which can be notated either $\text{div}(\text{grad } f)$ or $\nabla \cdot (\nabla f)$. The important thing to recognize is that the gradient, notated $\text{grad } f$ or ∇f , outputs a *vector*, not a scalar like the quantity Q defined in this problem.

Solutions for chapter 14

Page 1011, problem 2:

$\{\hat{x}\}$ is not a basis, because there are vectors such as \hat{y} that we can't form as a linear combination (i.e., scalar multiple) of \hat{x} . $\{\hat{x}, \hat{y}\}$ is the standard basis for this vector space. $\{-\hat{x}, \hat{x} + \hat{y}\}$ also

works as a basis, because the two vectors are linearly independent, and it's easy to check that any vector in the plane can be formed as a linear superposition of them. $\{\hat{x}, \hat{y}, \hat{x} + \hat{y}\}$ is not a basis, because these three vectors are not linearly independent.

Page 1011, problem 3:

- (a) The sketch for ℓ will be a 45-degree line through the origin, while r will be only the part of that line in the first quadrant. Of the two, only ℓ is a vector space. The set r isn't a vector space, because it doesn't have additive inverses.
- (b) We have $(1/2)(\pi + \pi) = 0$, but $(1/2)\pi + (1/2)\pi = \pi$.

Page 1011, problem 4:

To do anything useful with these expressions describing units, we need to be able to talk about things like dividing meters by seconds to get meters per second. Thus “addition” needs to be multiplication, which corresponds to adding the exponents. Scalar “multiplication” actually has to be exponentiation, e.g., “multiplying” units of meters by the scalar 2 should give square meters.

Page 1016, problem 20:

$$\begin{aligned}\langle \Psi | \mathcal{O}_E \Psi \rangle &= \langle c^* \curvearrowleft + c'^* \curvearrowright | c \curvearrowleft + 4c' \curvearrowright \rangle \\ &= c^* c \langle \curvearrowleft | \curvearrowleft \rangle + 4c^* c' \langle \curvearrowleft | \curvearrowright \rangle + 4c'^* c \langle \curvearrowright | \curvearrowleft \rangle + 4c'^* c' \langle \curvearrowright | \curvearrowright \rangle\end{aligned}$$

The second and third terms vanish, because \curvearrowleft and \curvearrowright are distinguishable. Because \curvearrowleft and \curvearrowright are normalized, the result reduces to

$$c^* c + 4c'^* c' = |c|^2 + 4|c'|^2 = \frac{1}{2}(1 + 4) = 2.5.$$

Answers to self-checks

Answers to self-checks for chapter 0

Page 15:

If only he has the special powers, then his results can never be reproduced.

Page 17:

They would have had to weigh the rays, or check for a loss of weight in the object from which they were have emitted. (For technical reasons, this was not a measurement they could actually do, hence the opportunity for disagreement.)

Page 25:

A dictionary might define “strong” as “possessing powerful muscles,” but that’s not an operational definition, because it doesn’t say how to measure strength numerically. One possible operational definition would be the number of pounds a person can bench press.

Page 27:

A microsecond is 1000 times longer than a nanosecond, so it would seem like 1000 seconds, or about 20 minutes.

Page 28:

Exponents have to do with multiplication, not addition. The first line should be 100 times longer than the second, not just twice as long.

Page 31:

The various estimates differ by 5 to 10 million. The CIA's estimate includes a ridiculous number of gratuitous significant figures. Does the CIA understand that every day, people are born in, die in, immigrate to, and emigrate from Nigeria?

Page 32:

- (1) 4; (2) 2; (3) 2

Page 35:

$$1 \text{ yd}^2 \times (3 \text{ ft}/1 \text{ yd})^2 = 9 \text{ ft}^2$$
$$1 \text{ yd}^3 \times (3 \text{ ft}/1 \text{ yd})^3 = 27 \text{ ft}^3$$

Page 41:

$$C_1/C_2 = (w_1/w_2)^4$$

Answers to self-checks for chapter 1**Page 58:**

The stream has to spread out. When the velocity becomes zero, it seems like the cross-sectional area has to become infinite. In reality, this is the point where the water turns around and comes back down. The infinity isn't real; it occurs mathematically because we used a simplified model of the stream of water, assuming, for instance, that the water's velocity is always straight up.

Page 60:

A positive Δx means the object is moving in the same direction as the positive x axis. A negative Δx means it's going the opposite direction.

Page 66:

(1) The effect only occurs during blastoff, when their velocity is changing. Once the rocket engines stop firing, their velocity stops changing, and they no longer feel any effect. (2) It is only an observable effect of your motion relative to the air.

Page 68:

Galilean relativity says that experiments can't come out differently just because they're performed while in motion. The tilting of the surface tells us the train is accelerating, but it doesn't tell us anything about the train's velocity at that instant. The person in the train might say the bottle's velocity was zero (but changing), whereas a person working in a reference frame attached to the dirt outside says it's moving; they don't agree on velocities. They *do* agree on accelerations. The person in the train has to agree that the train is accelerating, since otherwise there's no reason for the funny tilting effect.

Page 69:

Yes. In U.S. currency, for instance, the quantum of money is one cent.

Answers to self-checks for chapter 2**Page 90:**

There are two reasonable possibilities we could imagine — neither of which ends up making much sense — if we insist on the straight-line trajectory. (1) If the car has constant speed along the line, then in the * frame we see it going straight down at constant speed. It makes sense that it goes straight down in the * frame of reference, since in that frame it was never moving horizontally, and there's no reason for it to start. However, it doesn't make sense that it goes down with constant speed, since falling objects are supposed to speed up the whole time they

fall. This violates both Galilean relativity and conservation of energy. (2) If it's speeding up and moving along a diagonal line in the original frame, then it might be conserving energy in one frame or the other. But if it's speeding up along a line, then as seen in the original frame, both its vertical motion and its horizontal motion must be speeding up. If its horizontal velocity is increasing in the original frame, then it can't be zero and remain zero in the $*$ frame. This violates Galilean relativity, since in the $*$ frame the car apparently starts moving sideways for no reason.

Answers to self-checks for chapter 3

Page 146:

By shifting his weight around, he can cause the center of mass not to coincide with the geometric center of the wheel.

Page 155:

(1) This is motion, not force. (2) This is a description of how the sub is able to get the water to produce a forward force on it. (3) The sub runs out of energy, not force.

Page 156:

Frictionless (or nearly frictionless) ice can certainly make a normal force, since otherwise a hockey puck would sink into the ice. Friction is not possible without a normal force, however: we can see this from the equation, or from common sense, e.g. while sliding down a rope you don't get any friction unless you grip the rope.

Page 157:

(1) It's kinetic friction, because her uniform is sliding over the dirt. (2) It's static friction, because even though the two surfaces are moving relative to the landscape, they're not slipping over each other. (3) Only kinetic friction creates heat, as when you rub your hands together. If you move your hands up and down together without sliding them across each other, no heat is produced by the static friction.

Page 158:

(1) Normal forces are always perpendicular to the surface of contact, which means right or left in this figure. Normal forces are repulsive, so the cliff's force on the feet is to the right, i.e., away from the cliff. (2) Frictional forces are always parallel to the surface of contact, which means right or left in this figure. Static frictional forces are in the direction that would tend to keep the surfaces from slipping over each other. If the wheel was going to slip, its surface would be moving to the left, so the static frictional force on the wheel must be in the direction that would prevent this, i.e., to the right. This makes sense, because it is the static frictional force that accelerates the dragster. (3) Normal forces are always perpendicular to the surface of contact. In this diagram, that means either up and to the left or down and to the right. Normal forces are repulsive, so the ball is pushing the bat away from itself. Therefore the ball's force is down and to the right on this diagram.

Page 175:

The dashed lines on the graph are about twice as far apart in the second cycle compared to the first, so the amplitude has doubled. For sufficiently small oscillations around an equilibrium with $x = 0$ and $U(0) = 0$, it's always a good approximation to take $U \propto x^2$, so the energy is proportional to the square of the amplitude; this is a general fact about all oscillations, provided that the amplitude is small. Since the amplitude doubled, the energy quadrupled.

Page 178:

The two graphs start off with the same amplitude, but the solid curve loses amplitude more rapidly. For a given time, t , the quantity e^{-ct} is apparently smaller for the solid curve, meaning that ct is greater. The solid curve has the higher value of c .

Page 180:

A decaying exponential never dies out to zero in any finite amount of time.

Page 184:

In the expression

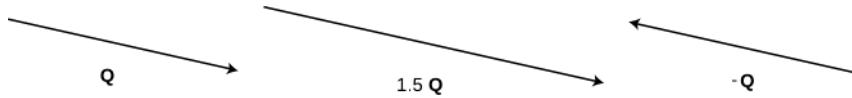
$$A = \frac{F_m}{m\sqrt{(\omega^2 - \omega_0^2)^2 + \omega_0^2\omega^2Q^{-2}}}$$

from page 1027, substituting $\omega = \omega_0$ makes the first term inside the square root vanish, which should make the denominator pretty small, thereby producing a pretty big amplitude. In the limit of $Q = \infty$, $Q^{-2} = 0$, so the second term vanishes, and $\omega = \omega_0$ actually produces an infinite amplitude. For values of Q that are large but finite, we still expect to get resonance pretty close to $\omega = \omega_0$. Setting $\omega = \omega_0$ in the finite- Q case, the first term vanishes, we can simplify the square root, and the result ends up being $A \propto 1/\sqrt{Q^{-2}} \propto Q$. This is only an approximation, because we had to assume early on that Q was large.

Page 198:

$$\mathbf{F} = m\mathbf{a}$$

Page 199:



Answers to self-checks for chapter 4

Page 262:

Torques 1, 2, and 4 all have the same sign, because they are trying to twist the wrench clockwise. The sign of 3 is opposite to the signs of 1, 2, and 4. The magnitude of 3 is the greatest, since it has a large r and the force is nearly all perpendicular to the wrench. Torques 1 and 2 are the same because they have the same values of r and F_{\perp} . Torque 4 is the smallest, due to its small r .

Page 271:

One person's θ - t graph would simply be shifted up or down relative to the others. The derivative equals the slope of the tangent line, and this slope isn't changed when you shift the graph, so both people would agree on the angular velocity.

Page 273:

Reversing the direction of ω also reverses the direction of motion, and this is reflected by the relationship between the plus and minus signs. In the equation for the radial acceleration, ω is squared, so even if ω is negative, the result is positive. This makes sense because the acceleration is always inward in circular motion, never outward.

Page 285:

All the rotations around the x axis give ω vectors along the positive x axis (thumb pointing along positive x), and all the rotations about the y axis have ω vectors with positive y components.

Page 288:

For example, if we take $(\mathbf{A} \times \mathbf{B})_x = A_y B_z - B_y A_z$ and reverse the A's and B's, we get $(\mathbf{B} \times \mathbf{A})_x = B_y A_z - A_y B_z$, which is just the negative of the original expression.

Answers to self-checks for chapter 5

Page 311:

Solids can exert shear forces. A solid could be in an equilibrium in which the shear forces were canceling the forces due to unequal pressures on the sides of the cube.

Page 311:

(1) Not valid. The equation only applies to fluids. (2) Valid. The density of the air is nearly constant between the top and bottom of the building. (3) Not valid. There is a large difference in the density of the air between the top and the bottom of the mountain. (4) Not valid, because g isn't constant throughout the interior of the earth. (5) Not valid, because the air is flowing around the wing. The air is accelerating, so it is not in equilibrium.

Page 331:

Heating the gas at constant pressure requires adding heat to it, which increases its entropy. To increase the gas's pressure while keeping its temperature constant, we would have to compress it, which would give it a smaller volume to inhabit, and therefore fewer possible positions for each atom. The whole thing has to be proportional to n because entropy is additive.

Answers to self-checks for chapter 6

Page 356:

The leading edge is moving up, the trailing edge is moving down, and the top of the hump is motionless for one instant.

Page 375:

The energy of a wave is usually proportional to the square of the amplitude. Squaring a negative number gives a positive result, so the energy is the same

Page 376:

A substance is invisible to sonar if the speed of sound waves in it is the same as in water. Reflections occur only at boundaries between media in which the wave speed is different.

Page 377:

No. A material object that loses kinetic energy slows down, but a wave is not a material object. The velocity of a wave ordinarily only depends on the medium, not on the amplitude. The speed of soft sound, for example, is the same as the speed of loud sound.

Page 383:

No. To get the best possible interference, the thickness of the coating must be such that the second reflected wave train lags behind the first by an integer number of wavelengths. Optimal performance can therefore only be produced for one specific color of light. The typical greenish color of the coatings shows that it does the worst job for green light.

Page 385:

The period is the time required to travel a distance $2L$ at speed v , $T = 2L/v$. The frequency is $f = 1/T = v/2L$.

Page 388:

The wave pattern will look like this:  Three quarters of a wavelength fit in the tube, so the wavelength is three times shorter than that of the lowest-frequency mode, in which one quarter of a wave fits. Since the wavelength is smaller by a factor of three, the frequency is three times

higher. Instead of f_o , $2f_o$, $3f_o$, $4f_o$, ..., the pattern of wave frequencies of this air column goes f_o , $3f_o$, $5f_o$, $7f_o$, ...

Answers to self-checks for chapter 7

Page 406:

At $v = 0$, we get $\gamma = 1$, so $t = T$. There is no time distortion unless the two frames of reference are in relative motion.

Page 428:

Both the time axis and the position axis have been turned around. Flipping the time axis means that the roles of transmitter and receiver have been swapped, and it also means that Alice and Betty are approaching one another rather than receding. The time experienced by the receiving observer is now the longer one, so the Doppler-shift factor has been inverted: the receiver now measures a Doppler shift of $1/2$ rather than 2 in frequency.

Page 430:

The total momentum is zero before the collision. After the collision, the two momenta have reversed their directions, but they still cancel. Neither object has changed its kinetic energy, so the total energy before and after the collision is also the same.

Page 437:

At $v = 0$, we have $\gamma = 1$, so the mass-energy is mc^2 as claimed. As v approaches c , γ approaches infinity, so the mass energy becomes infinite as well.

Answers to self-checks for chapter 8

Page 477:

Either type can be involved in either an attraction or a repulsion. A positive charge could be involved in either an attraction (with a negative charge) or a repulsion (with another positive), and a negative could participate in either an attraction (with a positive) or a repulsion (with a negative).

Page 478:

It wouldn't make any difference. The roles of the positive and negative charges in the paper would be reversed, but there would still be a net attraction.

Page 488:

Yes. In U.S. currency, the quantum of money is the penny.

Page 508:

Thomson was accelerating electrons, which are negatively charged. This apparatus is supposed to accelerate atoms with one electron stripped off, which have positive net charge. In both cases, a particle that is between the plates should be attracted by the forward plate and repelled by the plate behind it.

Page 517:

The hydrogen-1 nucleus is simple a proton. The binding energy is the energy required to tear a nucleus apart, but for a nucleus this simple there is nothing to tear apart.

Answers to self-checks for chapter 9

Page 547:

The large amount of power means a high rate of conversion of the battery's chemical energy into heat. The battery will quickly use up all its energy, i.e., "burn out."

Answers to self-checks for chapter 10

Page 585:

The reasoning is exactly analogous to that used in example 1 on page 583 to derive an equation for the gravitational field of the earth. The field is $F/q_t = (kQq_t/r^2)/q_t = kQ/r^2$.

Page 594:

$$\begin{aligned} E_x &= -\frac{dV}{dx} \\ &= -\frac{d}{dx} \left(\frac{kQ}{r} \right) \\ &= \frac{kQ}{r^2} \end{aligned}$$

Page 596:

- (a) The voltage (height) increases as you move to the east or north. If we let the positive x direction be east, and choose positive y to be north, then dV/dx and dV/dy are both positive. This means that E_x and E_y are both negative, which makes sense, since the water is flowing in the negative x and y directions (south and west).
- (b) The electric fields are all pointing away from the higher ground. If this was an electrical map, there would have to be a large concentration of charge all along the top of the ridge, and especially at the mountain peak near the south end.

Page 608:

- (a) The energy density depends on $\mathbf{E} \cdot \mathbf{E}$, which equals $E_x^2 + E_y^2 + E_z^2$.
- (b) Since E_x is squared, reversing its sign has no effect on the energy density. This makes sense, because otherwise we'd be saying that the positive and negative x axes in space were somehow physically different in their behavior, which would violate the symmetry of space.

Page 608:

$$\begin{aligned} N^{-1}m^{-2}C^2V^2m^{-2}m^2m &= N^{-1}m^{-1}C^2V^2 \\ &= N^{-1}m^{-1}J^2 \\ &= J^{-1}J^2 \\ &= J \end{aligned}$$

Page 617:

Yes. The mass has the same kinetic energy regardless of which direction it's moving. Friction converts mechanical energy into heat at the same rate whether the mass is sliding to the right or to the left. The spring has an equilibrium length, and energy can be stored in it either by compressing it ($x < 0$) or stretching it ($x > 0$).

Page 618:

Velocity, v , is the rate of change of position, x , with respect to time. This is exactly analogous to $I = \Delta q/\Delta t$.

Page 628:

Say we're looking for $u = \sqrt{z}$, i.e., we want a number u that, multiplied by itself, equals z .

Multiplication multiplies the magnitudes, so the magnitude of u can be found by taking the square root of the magnitude of z . Since multiplication also adds the arguments of the numbers, squaring a number doubles its argument. Therefore we can simply divide the argument of z by two to find the argument of u . This results in one of the square roots of z . There is another one, which is $-u$, since $(-u)^2$ is the same as u^2 . This may seem a little odd: if u was chosen so that doubling its argument gave the argument of z , then how can the same be true for $-u$? Well for example, suppose the argument of z is 4° . Then $\arg u = 2^\circ$, and $\arg(-u) = 182^\circ$. Doubling 182 gives 364, which is actually a synonym for 4 degrees.

Page 631:

Only $\cos(6t - 4)$ can be represented by a complex number. Although the graph of $\cos^2 t$ does have a sinusoidal shape, it varies between 0 and 1, rather than -1 and 1 , and there is no way to represent that using complex numbers. The function $\tan t$ doesn't even have a sinusoidal shape.

Page 632:

The impedance depends on the frequency at which the capacitor is being driven. It isn't just a single value for a particular capacitor.

Page 645:

The quantity $4\pi kq_{in}$ is now negative, so we'd better get a negative flux on the other side of Gauss' theorem. We do, because each field vector \mathbf{E}_j is inward, while the corresponding area vector, \mathbf{A}_j , is outward. Vectors in opposite directions make negative dot products.

Answers to self-checks for chapter 11

Page 684:

From the top panel of the figure, where the magnetic field is turned off, we can see that the beam leaves the cathode traveling upward, so in the bottom figure the electrons must be circling in the counterclockwise direction. To produce circular motion, the force must be towards the center of the circle. We can arbitrarily pick a point on the circle at which to analyze the vectors — let's pick the right-hand side. At this point, the velocity vector of the electrons is upward. Since the electrons are negatively charged, the force $q\mathbf{v} \times \mathbf{B}$ is given by $-\mathbf{v} \times \mathbf{B}$, not $+\mathbf{v} \times \mathbf{B}$. Circular orbits are produced when the motion is in the plane perpendicular to the field, so the field must be either into or out of the page. If the field was into the page, the right-hand rule would give $\mathbf{v} \times \mathbf{B}$ to the left, which is towards the center, but the force would be in the direction of $-\mathbf{v} \times \mathbf{B}$, which would be outwards. The field must be out of the page.

Page 686:

For instance, imagine a small sphere around the negative charge, which we would sketch on the two-dimensional paper as a circle. The field points inward at every point on the sphere, so all the contributions to the flux are negative. There is no cancellation, and the total flux is negative, which is consistent with Gauss' law, since the sphere encloses a negative charge. Copying the same surface onto the field of the bar magnet, however, we find that there is inward flux on the top and outward flux on the bottom, where the surface is inside the magnet. According to Gauss' law for magnetism, these cancel exactly, which is plausible based on the figure.

Page 690:

Plugging $z = 0$ into the equation gives $B_z = 4kI/c^2h$. This is simply twice the field of a single wire at a distance h . At this location, the fields contributed by the two wires are parallel, so vector addition simply gives a vector twice as strong.

Page 704:

The circulation around the Ampèrean surface we used was counterclockwise, since the field on the bottom was to the right. Applying the right-hand rule, the current I_{through} must have been out of the page at the top of the solenoid, and into the page at the bottom.

Page 704:

The quantity ℓ came in because we set $\eta = NI/\ell$. Based on that, it's clear that ℓ represents the length of the solenoid, not the length of the wire.

Page 704:

Doubling the radius of the solenoid would mean that every distance in the problem would be doubled, which would tend to make the fields weaker, since fields fall off with distance. However, doubling the radius would also mean that we had twice as much wire, and therefore twice as many moving charges to create magnetic fields. Since the magnetic field of a wire falls off like $1/r$, it's not surprising that the first effect amounts to exactly a factor of $1/2$, which is exactly enough to cancel out the factor of 2 from the second effect.

Page 716:

An induced electric field can only be created by a *changing* magnetic field. Nothing is changing if your car is just sitting there. A point on the coil won't experience a changing magnetic field unless the coil is already spinning, i.e., the engine has already turned over.

Page 721:

Let's get all the electrical units in terms of Teslas. Electric field units can be expressed as $T \cdot m/s$. The circulation of the electric field has units of electric field multiplied by distance, or $T \cdot m^2/s$. On the right side, the derivative $\partial \mathbf{B} / \partial t$ has units of T/s , and multiplying this by area gives units of $T \cdot m^2/s$, just like on the left side.

Page 739:

An (idealized) battery is a circuit element that always maintains the same voltage difference across itself, so by the loop rule, the voltage difference across the capacitor must remain unchanged, even while the dielectric is being withdrawn. The bound charges on the surfaces of the dielectric have been attracting the free charges in the plates, causing them to charge up more than they ordinarily would have. As the dielectric is withdrawn, the capacitor will be partially discharged, and we will observe a current in the ammeter. Since the dielectric is attracted to the plates, positive work is done in extracting it, indicating that there must be an increase in the electrical energy stored in the capacitor. This may seem paradoxical, since the energy stored in a capacitor is $(1/2)CV^2$, and we are decreasing the capacitance. However, the energy $(1/2)CV^2$ is calculated in terms of the work required to deposit the free charge on the plates. In addition to this energy, there is also energy stored in the dielectric itself. By moving its bound charges farther away from the free charges in the plates, to which they are attracted, we have increased their electrical energy. This energy of the bound charges is inaccessible to the electric circuit.

Answers to self-checks for chapter 12

Page 776:

Only 1 is correct. If you draw the normal that bisects the solid ray, it also bisects the dashed ray.

Page 780:

You should have found from your ray diagram that an image is still formed, and it has simply moved down the same distance as the real face. However, this new image would only be visible from high up, and the person can no longer see his own image.

Page 785:

Increasing the distance from the face to the mirror has decreased the distance from the image to the mirror. This is the opposite of what happened with the virtual image.

Page 795:

At the top of the graph, d_i approaches infinity when d_o approaches f . Interpretation: the rays just barely converge to the right of the mirror.

On the far right, d_i approaches f as d_o approaches infinity; this is the definition of the focal length.

At the bottom, d_i approaches negative infinity when d_o approaches f from the other side. Interpretation: the rays don't quite converge on the right side of the mirror, so they appear to have come from a virtual image point very far to the left of the mirror.

Page 804:

(1) If n_1 and n_2 are equal, Snell's law becomes $\sin \theta_1 = \sin \theta_2$, which implies $\theta_1 = \theta_2$, since both angles are between 0 and 90°. The graph would be a straight line along the diagonal of the graph. (2) The graph is farthest from the diagonal when the angles are large, i.e., when the ray strikes the interface at a grazing angle.

Page 809:

(1) In 1, the rays cross the image, so it's real. In 2, the rays only appear to have come from the image point, so the image is virtual. (2) A rays is always closer to the normal in the medium with the higher index of refraction. The first left turn makes the ray closer to the normal, which is what should happen in glass. The second left turn makes the ray farther from the normal, and that's what should happen in air. (3) Take the topmost ray as an example. It will still take two right turns, but since it's entering the lens at a steeper angle, it will also leave at a steeper angle. Tracing backward to the image, the steeper lines will meet closer to the lens.

Page 817:

It would have to have a wavelength on the order of centimeters or meters, the same distance scale as that of your body. These would be microwaves or radio waves. (This effect can easily be noticed when a person affects a TV's reception by standing near the antenna.) None of this contradicts the correspondence principle, which only states that the wave model must agree with the ray model when the ray model is applicable. The ray model is not applicable here because λ/d is on the order of 1.

Page 819:

At this point, both waves would have traveled nine and a half wavelengths. They would both be at a negative extreme, so there would be constructive interference.

Page 823:

Judging by the distance from one bright wave crest to the next, the wavelength appears to be about 2/3 or 3/4 as great as the width of the slit.

Page 824:

Since the wavelengths of radio waves are thousands of times longer, diffraction causes the resolution of a radio telescope to be thousands of times worse, all other things being equal. (To compensate for the wavelength, it's desirable to make the telescope very large, as in figure z on page 824.)

(1 rectangle = $5 \text{ cm} \times 0.005 \text{ cm}^{-1} = 0.025$), but that would have been pointless, because we were just going to compare the two areas.

Answers to self-checks for chapter 13

Page 862:

(1) Most people would think they were positively correlated, but it's possible that they're independent. (2) These must be independent, since there is no possible physical mechanism that could make one have any effect on the other. (3) These cannot be independent, since dying today guarantees that you won't die tomorrow.

Page 864:

The area under the curve from 130 to 135 cm is about 3/4 of a rectangle. The area from 135 to 140 cm is about 1.5 rectangles. The number of people in the second range is about twice as much. We could have converted these to actual probabilities (1 rectangle = $5 \text{ cm} \times 0.005 \text{ cm}^{-1} = 0.025$), but that would have been pointless because we were just going to compare the two areas.

Page 869:

On the left-hand side, dN is a unitless count, and dt is an infinitesimal amount of time, with units of seconds, so the units are s^{-1} as claimed. On the right, both $N(0)$ and the exponential factor are unitless, so the only units come from the factor of $1/\tau$, which again has units of s^{-1} .

Page 877:

The axes of the graph are frequency and photon energy, so its slope is Planck's constant. It doesn't matter if you graph $e\Delta V$ rather than $W + e\Delta V$, because that only changes the y-intercept, not the slope.

Page 895:

Wavelength is inversely proportional to momentum, so to produce a large wavelength we would need to use electrons with very small momenta and energies. (In practical terms, this isn't very easy to do, since ripping an electron out of an object is a violent process, and it's not so easy to calm the electrons down afterward.)

Page 904:

Under the ordinary circumstances of life, the accuracy with which we can measure position and momentum of an object doesn't result in a value of $\Delta p\Delta x$ that is anywhere near the tiny order of magnitude of Planck's constant. We run up against the ordinary limitations on the accuracy of our measuring techniques long before the uncertainty principle becomes an issue.

Page 907:

No. The equation $KE = p^2/2m$ is nonrelativistic, so it can't be applied to an electron moving at relativistic speeds. Photons always move at relativistic speeds, so it can't be applied to them, either.

Page 909:

Dividing by Planck's constant, a small number, gives a large negative result inside the exponential, so the probability will be very small.

Page 922:

If you trace a circle going around the center, you run into a series of eight complete wavelengths. Its angular momentum is $8\hbar$.

Page 928:

$n = 3, \ell = 0, \ell_z = 0$: one state; $n = 3, \ell = 1, \ell_z = -1, 0, \text{ or } 1$: three states; $n = 3, \ell = 2, \ell_z = -2, -1, 0, 1, \text{ or } 2$: five states

Page 934:

The original argument was that a kink would have a zero wavelength, which would correspond to an infinite momentum and an infinite kinetic energy, and that would violate conservation of energy. But the kink in this example occurs at $r = 0$, which is right on top of the proton, where the electrical energy $-ke^2/r$ is infinite and *negative*. Since the electrical energy is negative and infinite, we're actually *required* to have an infinite positive kinetic energy in order to come up with a total that conserves energy.

Answers

Answers for chapter 2

Page 126, problem 37:

$$K = k_1 k_2 / (k_1 + k_2) = 1 / (1/k_1 + 1/k_2)$$

Answers for chapter 3

Page 222, problem 5:

After the collision it is moving at $1/3$ of its initial speed, in the same direction it was initially going (it “follows through”).

Page 229, problem 41:

$$Q = 1/\sqrt{2}$$

Page 229, problem 43:

(a) 7×10^{-8} radians, or about 4×10^{-6} degrees.

Page 231, problem 51:

(a) $R = (2v^2/g) \sin \theta \cos \theta$ (c) 45°

Page 231, problem 51:

(a) $R = (2v^2/g) \sin \theta \cos \theta$ (c) 45°

Page 231, problem 52:

(a) The optimal angle is about 40° , and the resulting range is about 124 meters, which is about the length of a home run. (b) It goes about 9 meters farther. For comparison with reality, the stadium’s web site claims a home run goes about 11 meters farther there than in a sea-level stadium.

Page 231, problem 52:

(a) The optimal angle is about 40° , and the resulting range is about 124 meters, which is about the length of a home run. (b) It goes about 9 meters farther. For comparison with reality, the stadium’s web site claims a home run goes about 11 meters farther there than in a sea-level stadium.

Answers for chapter 5

Page 348, problem 9:

(a) $\sim 2 - 10\%$ (b) 5% (c) The high end for the body’s actual efficiency is higher than the limit imposed by the laws of thermodynamics. However, the high end of the 1-5 watt range quoted in the problem probably includes large people who aren’t just lying around. Still, it’s impressive that the human body comes so close to the thermodynamic limit.

Page 349, problem 10:

(a) Looking up the relevant density for air, and converting everything to mks, we get a frequency of 730 Hz. This is on the right order of magnitude, which is promising, considering the crudeness

of the approximation. (b) This brings the result down to 400 Hz, which is amazingly close to the observed frequency of 300 Hz.

Answers for chapter 6

Page 393, problem 8:

(a) $T = \mu\omega^2 r^2$

Page 393, problem 9:

(b) $g/2$

Page 394, problem 13:

Check: The actual length of a flute is about 66 cm.

Page 395, problem 17:

(a) $f = 4\alpha/(1 + \alpha)^2$ (b) $v_2 = \sqrt{v_1 v_3}$

Page 395, problem 17:

(a) $f = 4\alpha/(1 + \alpha)^2$ (b) $v_2 = \sqrt{v_1 v_3}$

Answers for chapter 10

Page 662, problem 22:

(a) $E = 2k\lambda/R$.

Answers for chapter 11

Page 747, problem 5:

(a) $I = \lambda v$.

Page 748, problem 6:

(b) $2kI_1 I_2 L/c^2 R$.

Answers for chapter 12

Page 844, problem 61:

f/ϵ

Page 845, problem 62:

$P = (1/2)(n^2 - 1)$

Answers for chapter 13

Page 950, problem 40:

about 10^{-34}

Appendix 4: Useful Data

.0.1 Notation and terminology, compared with other books

Almost all the notation and terminology in *Simple Nature* is standard, but there are some cases where there is no universal standard, and a very few cases where I've intentionally deviated from a universal standard. The notation used by physicists is also different from that used by electrical and mechanical engineers; I use physics terminology and notation (notably $\sqrt{-1} = i$, not j , and “torque” rather than “moment”), but employ the SI system of units used in engineering, rather than the cgs units favored by some physicists.

Nonstandard terminology:

Potential energy is referred to in this book as *interaction energy*, or according to its type: *gravitational energy*, *electrical energy*, etc.

The potential, in an electrical context, is referred to as *voltage*, e.g. I say that $V = kq/r$ is the voltage surrounding a point charge.

Heat and thermal energy are both referred to as *heat*. This is in keeping with casual usage among scientists, but formal written usage dictates the use of “thermal energy” to mean the kinetic energy an object has because of its molecules’ random motion, while “heat” is the transfer of thermal energy.

Notation for which there is no universal standard:

Kinetic energy is written K . Standard notation is K , T , or KE .

Interaction energy is written U . Standard notation is U , V , or PE .

The unit vectors are $\hat{\mathbf{x}}, \hat{\mathbf{y}}, \hat{\mathbf{z}}$. Standard notation is either $\hat{\mathbf{x}}, \hat{\mathbf{y}}, \hat{\mathbf{z}}$ or $\hat{\mathbf{i}}, \hat{\mathbf{j}}, \hat{\mathbf{k}}$.

Distance from an axis in cylindrical coordinates is R . A more common notation in math books is ρ , but this would conflict with the standard physics notation for the charge density.

Vibrations do not have very well standardized terminology or notation. I use “frequency” to refer to both f and ω , depending on the context to make it clear which is meant. The frequency of free, damped oscillations is ω_f , which is only approximately the same as $\omega_0 = \sqrt{k/m}$. The full width at half-maximum of the resonance peak (on a plot of energy versus frequency) is $\Delta\omega$.

The coupling constants for electricity and magnetism are written as k and k/c^2 . This is standard notation, but it would be more common in SI calculations to see everything expressed in terms of $\epsilon_0 = 1/4\pi k$ and $\mu_0 = 4\pi k/c^2$. Numerically, we have $k = 8.99 \times 10^9 \text{ N}\cdot\text{m}^2/\text{C}^2$ and $k/c^2 = 10^{-7} \text{ N}/\text{A}^2$, the latter being an exact relation.

.0.2 Notation and units

quantity	unit	symbol
distance	meter, m	$x, \Delta x$
time	second, s	$t, \Delta t$
mass	kilogram, kg	m
density	kg/m^3	ρ
force	newton, 1 N = 1 $\text{kg}\cdot\text{m/s}^2$	F
velocity	m/s	v
acceleration	m/s^2	a
gravitational field	$\text{J/kg}\cdot\text{m}$ or m/s^2	g
energy	joule, J	E (also electric field)
momentum	$\text{kg}\cdot\text{m/s}$	p
angular momentum	$\text{kg}\cdot\text{m}^2/\text{s}$ or $\text{J}\cdot\text{s}$	L (also inductance)
power	watt, 1 W = 1 J/s	P (also pressure)
pressure	$1 \text{ Pa} = 1 \text{ N/m}^2$	P (also power)
temperature	K	T (also period)
period	s	T (also temperature)
wavelength	m	λ
frequency	s^{-1} or Hz	f
charge	coulomb, C	q
voltage	volt, 1 V = 1 J/C	V
current	ampere, 1 A = 1 C/s	I
resistance	ohm, 1 Ω = 1 V/A	R
capacitance	farad, 1 F = 1 C/V	C
inductance	henry, 1 H = 1 V·s/A	L (also angular momentum)
electric field	V/m or N/C	E (also energy)
magnetic field	tesla, 1 T = 1 N·s/C·m	B
focal length	m	f
magnification	unitless	M
index of refraction	unitless	n
electron wavefunction	$\text{m}^{-3/2}$	Ψ

.0.3 Fundamental constants

gravitational constant	$G = 6.67 \times 10^{-11} \text{ J}\cdot\text{m/kg}^2$
Boltzmann constant	$k = 1.38 \times 10^{-23} \text{ J/K}$
Coulomb constant	$k = 8.99 \times 10^9 \text{ J}\cdot\text{m/C}^2$ or $\text{N}\cdot\text{m}^2/\text{C}^2$
quantum of charge	$e = 1.60 \times 10^{-19} \text{ C}$
speed of light	$c = 3.00 \times 10^8 \text{ m/s}$
Planck's constant	$h = 6.63 \times 10^{-34} \text{ J}\cdot\text{s}$

Note the use of the same notation, k , for both the Boltzmann constant and the Coulomb constant.

.0.4 Metric prefixes

M-	mega-	10^6
k-	kilo-	10^3
m-	milli-	10^{-3}
μ - (Greek mu)	micro-	10^{-6}
n-	nano-	10^{-9}
p-	pico-	10^{-12}
f-	femto-	10^{-15}

Note that the exponents go in steps of three. The exception is centi-, 10^{-2} , which is used only in the centimeter, and this doesn't require memorization, because a cent is 10^{-2} dollars.

.0.5 Nonmetric units

Nonmetric units in terms of metric ones:

1 inch	= 25.4 mm (by definition)
1 pound (lb)	= 4.5 newtons of force
1 scientific calorie	= 4.18 J
1 nutritional calorie	= 4.18×10^3 J
1 gallon	= 3.78×10^3 cm ³
1 horsepower	= 746 W

The pound is a unit of force, so it converts to newtons, not kilograms. A one-kilogram mass at the earth's surface experiences a gravitational force of $(1 \text{ kg})(9.8 \text{ m/s}^2) = 9.8 \text{ N} = 2.2 \text{ lb}$. The nutritional information on food packaging typically gives energies in units of calories, but those so-called calories are really kilocalories.

Relationships among U.S. units:

1 foot (ft)	= 12 inches
1 yard (yd)	= 3 feet
1 mile (mi)	= 5280 feet
1 ounce (oz)	= 1/16 pound

.0.6 The Greek alphabet

α	A	alpha	ι	I	iota	ρ	P	rho
β	B	beta	κ	K	kappa	σ	Σ	sigma
γ	Γ	gamma	λ	Λ	lambda	τ	T	tau
δ	Δ	delta	μ	M	mu	ν	Y	upsilon
ϵ	E	epsilon	ν	N	nu	ϕ	Φ	phi
ζ	Z	zeta	ξ	Ξ	xi	χ	X	chi
η	H	eta	\circ	O	omicron	ψ	Ψ	psi
θ	Θ	theta	π	Π	pi	ω	Ω	omega

.0.7 Subatomic particles

particle	mass (kg)	charge	radius (fm)
electron	9.109×10^{-31}	$-e$	$\lesssim 0.01$
proton	1.673×10^{-27}	$+e$	~ 1.1
neutron	1.675×10^{-27}	0	~ 1.1
neutrino	$\sim 10^{-39}$ kg	?	?

The radii of protons and neutrons can only be given approximately, since they have fuzzy

surfaces. For comparison, a typical atom is about a million fm in radius.

.0.8 Earth, moon, and sun

body	mass (kg)	radius (km)	radius of orbit (km)
earth	5.97×10^{24}	6.4×10^3	1.49×10^8
moon	7.35×10^{22}	1.7×10^3	3.84×10^5
sun	1.99×10^{30}	7.0×10^5	—

.0.9 The periodic table

¹ H																			² He				
³ Li	⁴ Be																	⁵ B	⁶ C	⁷ N	⁸ O	⁹ F	¹⁰ Ne
¹¹ Na	¹² Mg																	¹³ Al	¹⁴ Si	¹⁵ P	¹⁶ S	¹⁷ Cl	¹⁸ Ar
¹⁹ K	²⁰ Ca	²¹ Sc	²² Ti	²³ V	²⁴ Cr	²⁵ Mn	²⁶ Fe	²⁷ Co	²⁸ Ni	²⁹ Cu	³⁰ Zn	³¹ Ga	³² Ge	³³ As	³⁴ Se	³⁵ Br	³⁶ Kr						
³⁷ Rb	³⁸ Sr	³⁹ Y	⁴⁰ Zr	⁴¹ Nb	⁴² Mo	⁴³ Tc	⁴⁴ Ru	⁴⁵ Rh	⁴⁶ Pd	⁴⁷ Ag	⁴⁸ Cd	⁴⁹ In	⁵⁰ Sn	⁵¹ Sb	⁵² Te	⁵³ I	⁵⁴ Xe						
⁵⁵ Cs	⁵⁶ Ba	⁵⁷ La	*	⁷² Hf	⁷³ Ta	⁷⁴ W	⁷⁵ Re	⁷⁶ Os	⁷⁷ Ir	⁷⁸ Pt	⁷⁹ Au	⁸⁰ Hg	⁸¹ Tl	⁸² Pb	⁸³ Bi	⁸⁴ Po	⁸⁵ At	⁸⁶ Rn					
⁸⁷ Fr	⁸⁸ Ra	⁸⁹ Ac	**	¹⁰⁴ Rf	¹⁰⁵ Ha	¹⁰⁶	¹⁰⁷	¹⁰⁸	¹⁰⁹	¹¹⁰	¹¹¹	¹¹²	¹¹³	¹¹⁴	¹¹⁵	¹¹⁶	¹¹⁷	¹¹⁸					
																			* ⁵⁸ Ce ⁵⁹ Pr ⁶⁰ Nd ⁶¹ Pm ⁶² Sm ⁶³ Eu ⁶⁴ Gd ⁶⁵ Tb ⁶⁶ Dy ⁶⁷ Ho ⁶⁸ Er ⁶⁹ Tm ⁷⁰ Yb ⁷¹ Lu				
																			** ⁹⁰ Th ⁹¹ Pa ⁹² U ⁹³ Np ⁹⁴ Pu ⁹⁵ Am ⁹⁶ Cm ⁹⁷ Bk ⁹⁸ Cf ⁹⁹ Es ¹⁰⁰ Fm ¹⁰¹¹⁰² No ¹⁰³ Lr				

.0.10 Atomic masses

These atomic masses are given in atomic mass units (u), where by definition the mass of an atom of the isotope carbon-12 equals 12 u. One atomic mass unit is the same as about 1.66×10^{-27} kg. Data are only given for naturally occurring elements.

Ag	107.9	Eu	152.0	Mo	95.9	Sc	45.0
Al	27.0	F	19.0	N	14.0	Se	79.0
Ar	39.9	Fe	55.8	Na	23.0	Si	28.1
As	74.9	Ga	69.7	Nb	92.9	Sn	118.7
Au	197.0	Gd	157.2	Nd	144.2	Sr	87.6
B	10.8	Ge	72.6	Ne	20.2	Ta	180.9
Ba	137.3	H	1.0	Ni	58.7	Tb	158.9
Be	9.0	He	4.0	O	16.0	Te	127.6
Bi	209.0	Hf	178.5	Os	190.2	Ti	47.9
Br	79.9	Hg	200.6	P	31.0	Tl	204.4
C	12.0	Ho	164.9	Pb	207.2	Tm	168.9
Ca	40.1	In	114.8	Pd	106.4	U	238
Ce	140.1	Ir	192.2	Pt	195.1	V	50.9
Cl	35.5	K	39.1	Pr	140.9	W	183.8
Co	58.9	Kr	83.8	Rb	85.5	Xe	131.3
Cr	52.0	La	138.9	Re	186.2	Y	88.9
Cs	132.9	Li	6.9	Rh	102.9	Yb	173.0
Cu	63.5	Lu	175.0	Ru	101.1	Zn	65.4
Dy	162.5	Mg	24.3	S	32.1	Zr	91.2
Er	167.3	Mn	54.9	Sb	121.8		

Appendix 5: Summary

Notation and units are summarized on page 1071.

Chapter 0, Introduction and Review, page 13

Physics is the use of the scientific method to study the behavior of light and matter. The scientific method requires a cycle of theory and experiment, theories with both predictive and explanatory value, and reproducible experiments.

The metric system is a simple, consistent framework for measurement built out of the meter, the kilogram, and the second plus a set of prefixes denoting powers of ten. The most systematic method for doing conversions is shown in the following example:

$$370 \text{ ms} \times \frac{10^{-3} \text{ s}}{1 \text{ ms}} = 0.37 \text{ s}$$

Mass is a measure of the amount of a substance. Mass can be defined gravitationally, by comparing an object to a standard mass on a double-pan balance, or in terms of inertia, by comparing the effect of a force on an object to the effect of the same force on a standard mass. The two definitions are found experimentally to be proportional to each other to a high degree of precision, so we usually refer simply to “mass,” without bothering to specify which type.

A force is that which can change the motion of an object. The metric unit of force is the Newton, defined as the force required to accelerate a standard 1-kg mass from rest to a speed of 1 m/s in 1 s.

Scientific notation means, for example, writing 3.2×10^5 rather than 320000.

Writing numbers with the correct number of significant figures correctly communicates how accurate they are. As a rule of thumb, the final result of a calculation is no more accurate than, and should have no more significant figures than, the least accurate piece of data.

Nature behaves differently on large and small scales. Galileo showed that this results fundamentally from the way area and volume scale. Area scales as the second power of length, $A \propto L^2$, while volume scales as length to the third power, $V \propto L^3$.

An order of magnitude estimate is one in which we do not attempt or expect an exact answer. The main reason why the uninitiated have trouble with order-of-magnitude estimates is that the human brain does not intuitively make accurate estimates of area and volume. Estimates of area and volume should be approached by first estimating linear dimensions, which one’s brain has a feel for.

Velocity, dx/dt , measures how fast an object is moving. Acceleration, d^2x/dt^2 , measures how quickly its velocity is changing. For motion with constant acceleration, we have these useful

relations:

$$\begin{aligned}a &= \frac{\Delta v}{\Delta t} \\x &= \frac{1}{2}at^2 + v_0t + x_0 \\v_f^2 &= v_0^2 + 2a\Delta x\end{aligned}$$

Chapter 1, Conservation of Mass, page 55

Conservation laws are the foundation of physics. A conservation law states that a certain quantity can be neither created nor destroyed; the total amount of it remains the same.

Mass is a conserved quantity in classical physics, i.e. physics before Einstein. This is plausible, since we know that matter is composed of subatomic particles; if the particles are neither created or destroyed, then it makes sense that the total mass will remain the same. There are two ways of defining mass.

Gravitational mass is defined by measuring the effect of gravity on a particular object, and comparing with some standard object, taking care to test both objects at a location where the strength of gravity is the same.

Inertial mass is defined by measuring how much a particular object resists a change in its state of motion. For instance, an object placed on the end of a spring will oscillate if the spring is initially compressed, and a more massive object will take longer to complete one oscillation.

Inertial and gravitational mass are equivalent: experiments show to a very high degree of precision that any two objects with the same inertial mass have the same gravitational mass as well.

The definition of inertial mass depends on a correct but counterintuitive assumption: that an object resists a change in its state of motion. Most people intuitively believe that motion has a natural tendency to slow down. This cannot be correct as a general statement, because “to slow down” is not a well-defined concept unless we specify what we are measuring motion relative to. This insight is credited to Galileo, and the general principle of *Galilean relativity* states that the laws of physics are the same in all inertial frames of reference. In other words, there is no way to distinguish a moving frame of reference from one that is at rest. To establish which frames of reference are inertial, we first must find one inertial frame in which objects appear to obey Galilean relativity. The surface of the earth is an inertial frame to a reasonably good approximation, and the frame of reference of the stars is an even better one. Once we have found one inertial frame of reference, any other frame is inertial which is moving in a straight line at constant velocity relative to the first one. For instance, if the surface of the earth is an approximately inertial frame, then a train traveling in a straight line at constant speed is also approximately an inertial frame.

The unit of mass is the kilogram, which, along with the meter and the second, forms the basis for the SI system of units (also known as the mks system). A fundamental skill in science is to know the definitions of the most common metric prefixes, which are summarized on page 1072, and to be able to convert among them.

One consequence of Einstein’s theory of special relativity is that *mass can be converted to energy and energy to mass*. This prediction has been verified amply by experiment. Thus the conserved quantity is not really mass but rather the total “mass-energy,” $m + E/c^2$, where c is the speed of light. Since the speed of light is a large number, the E/c^2 term is ordinarily small

in everyday life, which is why we can usually neglect it.

Chapter 2, Conservation of Energy, page 73

We observe that certain processes are physically impossible. For example, there is no process that can heat up an object without using up fuel or having some other side effect such as cooling a different object. We find that we can neatly separate the possible processes from the impossible by defining a single numerical quantity, called *energy*, which is conserved. Energy comes in many forms, such as heat, motion, sound, light, the energy required to melt a solid, and gravitational energy (e.g. the energy that depends on the distance between a rock and the earth). Because it has so many forms, we can arbitrarily choose one form, heat, in order to define a standard unit for our numerical scale of energy. Energy is measured in units of joules (J), and one joule can be defined as the amount of energy required in order to raise the temperature of a certain amount of water by a certain number of degrees. (The numbers are not worth memorizing.) *Power* is defined as the rate of change of energy $P = dE/dt$, and the unit of power is the watt, $1 \text{ W} = 1 \text{ J/s}$.

Once we have defined one type of energy numerically, we can perform experiments that establish the mathematical rules governing other types of energy. For example, in his paddlewheel experiment, James Joule allowed weights to drop through a certain height and spin paddlewheels inside sealed canisters of water, thereby heating the water through friction. Since in this book we define the joule unit in terms of the temperature of water, we can think of the paddlewheel experiment as establishing a rule for the *gravitational energy* of a mass which is at a certain height,

$$dU_g = mg dy,$$

where dU_g is the infinitesimal change in the gravitational energy of a mass m when its height is changed by an infinitesimal amount dy in the vertical direction. The quantity g is called the *gravitational field*, and at the earth's surface it has a numerical value of about $10 \text{ J/kg}\cdot\text{m}$. That is, about 10 joules of energy are required in order to raise a one-kilogram mass by one meter. (The gravitational field g also has the interpretation that when we drop an object, its acceleration, d^2y/dt^2 , is equal to g .)

Using similar techniques, we find that the energy of a moving object, called its *kinetic energy*, is given by

$$K = \frac{1}{2}mv^2,$$

where m is its mass and v its velocity. The proportionality factor equals $1/2$ exactly by the design of the SI system of units, and since the SI is based on the meter, the kilogram, and the second, the joule is considered to be a derived unit, $1 \text{ J} = 1 \text{ kg}\cdot\text{m}^2/\text{s}^2$.

When the interaction energy U has a local maximum or minimum with respect to the position of an object ($dU/dx = 0$), then the object is in *equilibrium* at that position. For example, if a weight is hanging from a rope, and is initially at rest at the bottom, then it must remain at rest, because this is a position of minimum gravitational energy U_g ; to move, it would have to increase both its kinetic and its gravitational energy, which would violate conservation of energy, since the total energy would increase.

Since kinetic energy is independent of the direction of motion, conservation of energy is often insufficient to predict the direction of an object's motion. However, many of the physically impossible motions can be ruled out by the trick of imposing conservation of energy in some other frame of reference. By this device, we can solve the important problem of *projectile motion*: even if the projectile has horizontal motion, we can imagine ourselves in a frame of reference in which

we are moving along with the projectile horizontally. In this frame of reference, the projectile has no horizontal motion, and its vertical motion has constant acceleration g . Switching back to a frame of reference in which its horizontal velocity is not zero, we find that a projectile's horizontal and vertical motions are independent, and that the horizontal motion is at constant velocity.

Even in one-dimensional motion, it is seldom possible to solve real-world problems and predict the motion of an object in closed form. However, there are straightforward numerical techniques for solving such problems.

From observations of the motion of the planets, we infer that the *gravitational interaction between any two objects* is given by $U_g = -Gm_1m_2/r$, where r is the distance between them. When the sizes of the objects are not small compared to their separation, the definition of r becomes vague; for this reason, we should interpret this fundamentally as the law governing the gravitational interactions between individual atoms. However, in the special case of a spherically symmetric mass distribution, there is a shortcut: the *shell theorem* states that the gravitational interaction between a spherically symmetric shell of mass and a particle on the outside of the shell is the same as if the shell's mass had all been concentrated at its center. An astronomical body like the earth can be broken down into concentric shells of mass, and so its gravitational interactions with external objects can also be calculated simply by using the center-to-center distance.

Energy appears to come in a bewildering variety of forms, but matter is made of atoms, and thus if we restrict ourselves to the study of mechanical systems (containing material objects, not light), all the forms of energy we observe must be explainable in terms of the behavior and interactions of atoms. Indeed, at the atomic level the picture is much simpler. Fundamentally, all the familiar forms of mechanical energy arise from either the kinetic energy of atoms or the energy they have because they interact with each other via gravitational or electrical forces. For example, when we stretch a spring, we distort the latticework of atoms in the metal, and this change in the interatomic distances involves an increase in the atoms' electrical energies.

An equilibrium is a local minimum of $U(x)$, and up close, any minimum looks like a parabola. Therefore, small oscillations around an equilibrium exhibit universal behavior, which depends only on the object's mass, m , and on the tightness of curvature of the minimum, parametrized by the quantity $k = d^2U/dx^2$. The oscillations are sinusoidal as a function of time, and the period is $T = 2\pi\sqrt{m/k}$, independent of amplitude. When oscillations are small enough for these statements to be good approximations, we refer to them the oscillations as *simple harmonic*.

Chapter 3, Conservation of Momentum, page 131

Since the kinetic energy of a material object depends on v^2 , it isn't obvious that conservation of energy is consistent with Galilean relativity. Even if a certain mechanical system conserves energy in one frame of reference, the velocities involved will be different as measured in another frame, and therefore so will the kinetic energies. It turns out that consistency is achieved only if there is a new conservation law, conservation of *momentum*,

$$\mathbf{p} = m\mathbf{v}.$$

In one dimension, the direction of motion is described using positive and negative signs of the velocity \mathbf{v} , and since mass is always positive, the momentum carries the same sign. Thus conservation of momentum, unlike conservation of energy, makes direct predictions about the direction of motion. Although this line of argument was based on the assumption of a mechanical system, momentum need not be mechanical. Light has momentum.

A moving object's momentum equals the sum of the momenta of all its atoms. To avoid having to carry out this sum, we can use the concept of the *center of mass*. The center of mass can be defined as a kind of weighted average of the positions of all the atoms in the object,

$$\mathbf{x}_{cm} = \frac{\sum m_j \mathbf{x}_j}{\sum m_j},$$

and although the definition does involve a sum, we can often locate the center of mass by symmetry or by physically determining an object's balance point. The total momentum of the object is then given by

$$\mathbf{P}_{total} = m_{total} \mathbf{v}_{cm}.$$

The rate of transfer of momentum is called *force*, $\mathbf{F} = d\mathbf{p}/dt$, and is measured in units of newtons, $1 \text{ N} = 1 \text{ kg}\cdot\text{m/s}^2$. As a direct consequence of conservation of momentum, we have the following statements, known as *Newton's laws of motion*:

If the total force on an object is zero, it remains in the same state of motion.

$$\mathbf{F} = d\mathbf{p}/dt$$

Forces always come in pairs: if object A exerts a force on object B, then object B exerts a force on object A which is the same strength, but in the opposite direction.

Although the fundamental forces at the atomic level are gravity, electromagnetism, and nuclear forces, we use a different and more practical classification scheme in everyday situations. In this scheme, the forces between solid objects are described as follows:

<i>A normal force</i> , F_n ,	is perpendicular to the surface of contact, and prevents objects from passing through each other by becoming as strong as necessary (up to the point where the objects break). “Normal” means perpendicular.
<i>Static friction</i> , F_s ,	is parallel to the surface of contact, and prevents the surfaces from starting to slip by becoming as strong as necessary, up to a maximum value of $F_{s,max}$. “Static” means not moving, i.e. not slipping.
<i>Kinetic friction</i> , F_k ,	is parallel to the surface of contact, and tends to slow down any slippage once it starts. “Kinetic” means moving, i.e. slipping.

Work is defined as the transfer of energy by a force. (“By a force” is meant to exclude energy transfer by heat conduction.) The *work theorem* states that when a force occurs at a single point of contact, the amount of energy transferred by that force is given by $dW = \mathbf{F} \cdot d\mathbf{x}$, where $d\mathbf{x}$ is the distance traveled by the point of contact. The *kinetic energy theorem* is $dK_{cm} = \mathbf{F}_{total} \cdot d\mathbf{x}_{cm}$, where dK_{cm} is the change in the energy, $(1/2)mv_{cm}^2$, an object possesses due to the motion of its center of mass, \mathbf{F}_{total} is the total force acting on the object, and $d\mathbf{x}_{cm}$ is the distance traveled by the center of mass.

The relationship between force and interaction energy is $U = -dF/dx$. Any interaction can be described either by giving the force as a function of distance or the interaction energy as a function of distance; the other quantity can then be found by integration or differentiation.

An oscillator subject to friction will, if left to itself, suffer a gradual decrease in the amplitude of its motion as mechanical energy is transformed into heat. The *quality factor*, Q , is defined as

the number of oscillations required for the mechanical energy to fall off by a factor of $e^{2\pi} \approx 535$. To maintain an oscillation indefinitely, an external force must do work to replace this energy. We assume for mathematical simplicity that the external force varies sinusoidally with time, $F = F_m \sin \omega t$. If this force is applied for a long time, the motion approaches a steady state, in which the oscillator's motion is sinusoidal, matching the driving force in frequency but not in phase. The amplitude of this steady-state motion, A , exhibits the phenomenon of *resonance*: the amplitude is maximized at a driving frequency which, for large Q , is essentially the same as the natural frequency of the free vibrations, ω_f (and for large Q this is also nearly the same as $\omega_0 = \sqrt{k/m}$). When the energy of the steady-state oscillations is graphed as a function of frequency, both the height and the width of the resonance peak depend on Q . The peak is taller for greater Q , and its full width at half-maximum is $\Delta\omega \approx \omega_0/Q$. For small values of Q , all these approximations become worse, and at $Q < 1/2$ qualitatively different behavior sets in.

For three-dimensional motion, a moving object's motion can be described by three different velocities, $v_x = dx/dt$, and similarly for v_y and v_z . Thus conservation of momentum becomes three different conservation laws: conservation of $p_x = mv_x$, and so on. The principle of *rotational invariance* says that the laws of physics are the same regardless of how we change the orientation of our laboratory: there is no preferred direction in space. As a consequence of this, no matter how we choose our x , y , and z coordinate axes, we will still have conservation of p_x , p_y , and p_z . To simplify notation, we define a momentum *vector*, \mathbf{p} , which is a single symbol that stands for all the momentum information contained in the components p_x , p_y , and p_z . The concept of a vector is more general than its application to the momentum: any quantity that has a direction in space is considered a vector, as opposed to a *scalar* like time or temperature. The following table summarizes some vector operations.

operation	definition
$ \mathbf{vector} $	$\sqrt{\text{vector}_x^2 + \text{vector}_y^2 + \text{vector}_z^2}$
$\mathbf{vector} + \mathbf{vector}$	Add component by component.
$\mathbf{vector} - \mathbf{vector}$	Subtract component by component.
$\mathbf{vector} \cdot \text{scalar}$	Multiply each component by the scalar.
$\mathbf{vector} / \text{scalar}$	Divide each component by the scalar.

Differentiation and integration of vectors is defined component by component.

There is only one meaningful (rotationally invariant) way of defining a multiplication of vectors whose result is a scalar, and it is known as the vector *dot product*:

$$\begin{aligned}\mathbf{b} \cdot \mathbf{c} &= b_x c_x + b_y c_y + b_z c_z \\ &= |\mathbf{b}| |\mathbf{c}| \cos \theta_{bc}.\end{aligned}$$

The dot product has most of the usual properties associated with multiplication, except that there is no “dot division.”

Chapter 4, Conservation of Angular Momentum, page 251

Angular momentum is a conserved quantity. For motion confined to a plane, the angular momentum of a material particle is

$$L = mv_{\perp}r,$$

where r is the particle's distance from the point chosen as the axis, and v_{\perp} is the component of its velocity vector that is perpendicular to the line connecting the particle to the axis. The choice of axis is arbitrary. In a plane, only two directions of rotation are possible, clockwise and counterclockwise. One of these is considered negative and the other positive. Geometrically,

angular momentum is related to rate at which area is swept out by the line segment connecting the particle to the axis.

Torque is the rate of change of angular momentum, $\tau = dL/dt$. The torque created by a given force can be calculated using any of the relations

$$\begin{aligned}\tau &= rF \sin \theta_{rF} \\ &= rF_{\perp} \\ &= r_{\perp} F,\end{aligned}$$

where the subscript \perp indicates a component perpendicular to the line connecting the axis to the point of application of the force.

In the special case of a *rigid body* rotating in a single plane, we define

$$\omega = \frac{d\theta}{dt} \quad [\text{angular velocity}]$$

and

$$\alpha = \frac{d\omega}{dt}, \quad [\text{angular acceleration}]$$

in terms of which we have

$$L = I\omega$$

and

$$\tau = I\alpha,$$

where the *moment of inertia*, I , is defined as

$$I = \sum m_i r_i^2,$$

summing over all the atoms in the object (or using calculus to perform a continuous sum, i.e. an integral). The relationship between the angular quantities and the linear ones is

$v_t = \omega r$	[tangential velocity of a point]
$v_r = 0$	[radial velocity of a point]
$a_t = \alpha r$.	[radial acceleration of a point]
	at a distance r from the axis]
$a_r = \omega^2 r$	[radial acceleration of a point]
	at a distance r from the axis]

In three dimensions, torque and angular momentum are vectors, and are expressed in terms of the vector *cross product*, which is the only rotationally invariant way of defining a multiplication of two vectors that produces a third vector:

$$\mathbf{L} = \mathbf{r} \times \mathbf{p}$$

$$\boldsymbol{\tau} = \mathbf{r} \times \mathbf{F}$$

In general, the cross product of vectors \mathbf{b} and \mathbf{c} has magnitude

$$|\mathbf{b} \times \mathbf{c}| = |\mathbf{b}| |\mathbf{c}| \sin \theta_{bc},$$

which can be interpreted geometrically as the area of the parallelogram formed by the two vectors when they are placed tail-to-tail. The direction of the cross product lies along the line which is perpendicular to both vectors; of the two such directions, we choose the one that is right-handed, in the sense that if we point the fingers of the flattened right hand along \mathbf{b} , then bend the knuckles to point the fingers along \mathbf{c} , the thumb gives the direction of $\mathbf{b} \times \mathbf{c}$. In terms of components, the cross product is

$$\begin{aligned} (\mathbf{b} \times \mathbf{c})_x &= b_y c_z - c_y b_z \\ (\mathbf{b} \times \mathbf{c})_y &= b_z c_x - c_z b_x \\ (\mathbf{b} \times \mathbf{c})_z &= b_x c_y - c_x b_y \end{aligned}$$

The cross product has the disconcerting properties

$$\mathbf{a} \times \mathbf{b} = -\mathbf{b} \times \mathbf{a} \quad [\text{noncommutative}]$$

and

$$\mathbf{a} \times (\mathbf{b} \times \mathbf{c}) \neq (\mathbf{a} \times \mathbf{b}) \times \mathbf{c} \quad [\text{nonassociative}],$$

and there is no “cross-division.”

For rigid-body rotation in three dimensions, we define an angular velocity vector $\boldsymbol{\omega}$, which lies along the axis of rotation and bears a right-hand relationship to it. Except in special cases, there is no scalar moment of inertia for which $\mathbf{L} = I\boldsymbol{\omega}$; the moment of inertia must be expressed as a matrix.

Chapter 5, Thermodynamics, page 307

A fluid is any gas or liquid, but not a solid; fluids do not exhibit shear forces. A fluid in equilibrium exerts a force on any surface which is proportional to the surface’s area and perpendicular to the surface. We can therefore define a quantity called the *pressure*, P , which is ratio of force to area,

$$P = \frac{F_\perp}{A},$$

where the subscript \perp indicates the component of the fluid’s force which is perpendicular to the surface.

Usually it is only the difference in pressure between the two sides of a surface that is physically significant. Pressure doesn’t just “press down” on things; air pressure upward under your chin is the same as air pressure downward on your shoulders. In a fluid acted on by gravity, pressure varies with depth according to the equation

$$dP = \rho \mathbf{g} \cdot dy.$$

This equation is only valid if the fluid is in equilibrium, and if g and r are constant with respect to height.

Temperature can be defined according to the volume of an ideal gas under conditions of standard pressure. The Kelvin scale of temperature used throughout this book equals zero at

absolute zero, the temperature at which all random molecular motion ceases, and equals 273 K at the freezing point of water. We can get away with using the Celsius scale as long as we are only interested in temperature differences; a difference of 1 degree C is the same as a difference of 1 degree K.

It is an observed fact that ideal gases obey the *ideal gas law*,

$$PV = nkT,$$

and this equation can be explained by the kinetic theory of heat, which states that the gas exerts pressure on its container because its molecules are constantly in motion. In the kinetic theory of heat, the temperature of a gas is proportional to the average energy per molecule.

Not all the heat energy in an object can be extracted to do mechanical work. We therefore describe heat as a lower grade of energy than other forms of energy. Entropy is a measure of how much of a system's energy is inaccessible to being extracted, even by the most efficient heat engine; a high entropy corresponds to a low grade of energy. The change in a system's entropy when heat Q is deposited into it is

$$\Delta S = \frac{Q}{T}.$$

The efficiency of any heat engine is defined as

$$\text{efficiency} = \frac{\text{energy we get in useful form}}{\text{energy we pay for}},$$

and the efficiency of a Carnot engine, the most efficient of all, is

$$\text{efficiency} = 1 - \frac{T_L}{T_H}.$$

These results are all closely related. For instance, example 11 on page 323 uses $\Delta S = Q/T$ and $\text{efficiency} = 1 - T_L/T_H$ to show that a Carnot engine doesn't change the entropy of the universe.

Fundamentally, entropy is defined as the being proportional to the natural logarithm of the number of states available to a system, and the above equation then serves as a definition of temperature. The entropy of a closed system always increases; this is the second law of thermodynamics.

Chapter 6, Waves, page 353

Wave motion differs in three important ways from the motion of material objects:

Waves obey the principle of superposition. When two waves collide, they simply add together.

The medium is not transported along with the wave. The motion of any given point in the medium is a vibration about its equilibrium location, not a steady forward motion.

The velocity of a wave depends on the medium, not on the amount of energy in the wave. (For some types of waves, notably water waves, the velocity may also depend on the shape of the wave.)

Sound waves consist of increases and decreases (typically very small ones) in the density of the air. Light is a wave, but it is a vibration of electric and magnetic fields, not of any physical medium. Light can travel through a vacuum.

A periodic wave is one that creates a periodic motion in a receiver as it passes it. Such a wave has a well-defined period and frequency, and it will also have a wavelength, which is the distance in space between repetitions of the wave pattern. The velocity, frequency, and wavelength of a periodic wave are related by the equation

$$v = f\lambda.$$

A wave emitted by a moving source will undergo a *Doppler shift* in wavelength and frequency. The shifted wavelength is given by the equation

$$\lambda' = \left(1 - \frac{u}{v}\right)\lambda,$$

where v is the velocity of the waves and u is the velocity of the source, taken to be positive or negative so as to produce a Doppler-lengthened wavelength if the source is receding and a Doppler-shortened one if it approaches. A similar shift occurs if the observer is moving, and in general the Doppler shift depends approximately only on the relative motion of the source and observer if their velocities are both small compared to the waves' velocity. (This is not just approximately but exactly true for light waves, as required by Einstein's theory of relativity.)

Whenever a wave encounters the boundary between two media in which its speeds are different, part of the wave is reflected and part is transmitted. The reflection is always reversed front-to-back, but may also be inverted in amplitude. Whether the reflection is inverted depends on how the wave speeds in the two media compare, e.g. a wave on a string is uninverted when it is reflected back into a segment of string where its speed is lower. The greater the difference in wave speed between the two media, the greater the fraction of the wave energy that is reflected. Surprisingly, a wave in a dense material like wood will be strongly reflected back into the wood at a wood-air boundary.

A one-dimensional wave confined by highly reflective boundaries on two sides will display motion which is periodic. For example, if both reflections are inverting, the wave will have a period equal to twice the time required to traverse the region, or to that time divided by an integer. An important special case is a sinusoidal wave; in this case, the wave forms a stationary pattern composed of a superposition of sine waves moving in opposite direction.

Chapter 7, Relativity, page 397

Experiments show that space and time do not have the properties claimed by Galileo and Newton. Time and space as seen by one observer are distorted compared to another observer's perceptions if they are moving relative to each other. This distortion is quantified by the factor

$$\gamma = \frac{1}{\sqrt{1 - \frac{v^2}{c^2}}},$$

where v is the relative velocity of the two observers, and c is a universal velocity that is the same in all frames of reference. Light travels at c . A clock appears to run fastest to an observer who is not in motion relative to it, and appears to run too slowly by a factor of γ to an observer who has a velocity v relative to the clock. Similarly, a meter-stick appears longest to an observer who sees it at rest, and appears shorter to other observers. Time and space are relative, not absolute. In particular, there is no well-defined concept of simultaneity.

All of these strange effects, however, are very small when the relative velocities are small compared to c . This makes sense, because Newton's laws have already been thoroughly tested by experiments at such speeds, so a new theory like relativity must agree with the old one in their realm of common applicability. This requirement of backwards-compatibility is known as the correspondence principle.

Relativity has implications not just for time and space but also for the objects that inhabit time and space. The correct relativistic equation for momentum is

$$p = m\gamma v,$$

which is similar to the classical $p = mv$ at low velocities, where $\gamma \approx 1$, but diverges from it more and more at velocities that approach c . Since γ becomes infinite at $v = c$, an infinite force would be required in order to give a material object enough momentum to move at the speed of light. In other words, material objects can only move at speeds lower than c . Relativistically, mass and energy are not separately conserved. Mass and energy are two aspects of the same phenomenon, known as mass-energy, and they can be converted to one another according to the equation

$$E = mc^2.$$

The mass-energy of a moving object is $E = m\gamma c^2$. When an object is at rest, $\gamma = 1$, and the mass-energy is simply the energy-equivalent of its mass, mc^2 . When an object is in motion, the excess mass-energy, in addition to the mc^2 , can be interpreted as its kinetic energy.

Chapter 8, Atoms and Electromagnetism, page 473

All the forces we encounter in everyday life boil down to two basic types: gravitational forces and electrical forces. A force such as friction or a “sticky force” arises from electrical forces between individual atoms.

Just as we use the word mass to describe how strongly an object participates in gravitational forces, we use the word *charge* for the intensity of its electrical forces. There are two types of charge. Two charges of the same type repel each other, but objects whose charges are different attract each other. Charge is measured in units of coulombs (C).

Mobile charged particle model: A great many phenomena are easily understood if we imagine matter as containing two types of charged particles, which are at least partially able to move around.

Positive and negative charge: Ordinary objects that have not been specially prepared have both types of charge spread evenly throughout them in equal amounts. The object will then tend not to exert electrical forces on any other object, since any attraction due to one type of charge will be balanced by an equal repulsion from the other. (We say “tend not to” because bringing the object near an object with unbalanced amounts of charge could cause its charges to separate from each other, and the force would no longer cancel due to the unequal distances.) It therefore makes sense to describe the two types of charge using positive and negative signs, so that an unprepared object will have zero *total* charge.

The *Coulomb force law* states that the magnitude of the electrical force between two charged particles is given by

$$|F| = \frac{k|q_1||q_2|}{r^2}.$$

Conservation of charge: An even more fundamental reason for using positive and negative

signs for charge is that with this definition the total charge of a closed system is a conserved quantity.

Quantization of charge: Millikan's oil drop experiment showed that the total charge of an object could only be an integer multiple of a basic unit of charge, e . This supported the idea that the "flow" of electrical charge was the motion of tiny particles rather than the motion of some sort of mysterious electrical fluid.

Einstein's analysis of Brownian motion was the first definitive proof of the existence of atoms. Thomson's experiments with vacuum tubes demonstrated the existence of a new type of microscopic particle with a very small ratio of mass to charge. Thomson correctly interpreted these as building blocks of matter even smaller than atoms: the first discovery of subatomic particles. These particles are called electrons.

The above experimental evidence led to the first useful model of the interior structure of atoms, called the raisin cookie model. In the raisin cookie model, an atom consists of a relatively large, massive, positively charged sphere with a certain number of negatively charged electrons embedded in it.

Rutherford and Marsden observed that some alpha particles from a beam striking a thin gold foil came back at angles up to 180 degrees. This could not be explained in the then-favored raisin-cookie model of the atom, and led to the adoption of the planetary model of the atom, in which the electrons orbit a tiny, positively-charged nucleus. Further experiments showed that the nucleus itself was a cluster of positively-charged protons and uncharged neutrons.

Radioactive nuclei are those that can release energy. The most common types of radioactivity are alpha decay (the emission of a helium nucleus), beta decay (the transformation of a neutron into a proton or vice-versa), and gamma decay (the emission of a type of very-high-frequency light). Stars are powered by nuclear fusion reactions, in which two light nuclei collide and form a bigger nucleus, with the release of energy.

Human exposure to ionizing radiation is measured in units of millirem. The typical person is exposed to about 100 mrem worth of natural background radiation per year.

Chapter 9, DC Circuits, page 531

All electrical phenomena are alike in that they arise from the presence or motion of charge. Most practical electrical devices are based on the motion of charge around a complete circuit, so that the charge can be recycled and does not hit any dead ends. The most useful measure of the flow of charge is *current*,

$$I = \frac{dq}{dt}.$$

An electrical device whose job is to transform energy from one form into another, e.g. a lightbulb, uses power at a rate which depends both on how rapidly charge is flowing through it and on how much work is done on each unit of charge. The latter quantity is known as the voltage difference between the point where the current enters the device and the point where the current leaves it. Since there is a type of electrical energy associated with electrical forces, the amount of work they do is equal to the difference in potential energy between the two points, and we therefore define voltage differences directly in terms of electrical energy,

$$\Delta V = \frac{\Delta U_{elec}}{q}.$$

The rate of power dissipation is

$$P = I\Delta V.$$

Many important electrical phenomena can only be explained if we understand the mechanisms of current flow at the atomic level. In metals, currents are carried by electrons, in liquids by ions. Gases are normally poor conductors unless their atoms are subjected to such intense electrical forces that the atoms become ionized.

Many substances, including all solids, respond to electrical forces in such a way that the flow of current between two points is proportional to the voltage difference between those points (assuming the voltage difference is small). Such a substance is called ohmic, and an object made out of an ohmic substance can be rated in terms of its resistance,

$$R = \frac{\Delta V}{I}$$

An important corollary is that a perfect conductor, with $R = 0$, must have constant voltage everywhere within it.

A schematic is a drawing of a circuit that standardizes and stylizes its features to make it easier to understand. Any circuit can be broken down into smaller parts. For instance, one big circuit may be understood as two small circuits in series, another as three circuits in parallel. When circuit elements are combined in parallel and in series, we have two basic rules to guide us in understanding how the parts function as a whole:

The junction rule: In any circuit that is not storing or releasing charge, conservation of charge implies that the total current flowing out of any junction must be the same as the total flowing in.

The loop rule: Assuming the standard convention for plus and minus signs, the sum of the voltage drops around any closed loop in a circuit must be zero.

The simplest application of these rules is to pairs of resistors combined in series or parallel. In such cases, the pair of resistors acts just like a single unit with a certain resistance value, called their equivalent resistance. Resistances in series add to produce a larger equivalent resistance,

$$R = R_1 + R_2,$$

because the current has to fight its way through both resistances. Parallel resistors combine to produce an equivalent resistance that is smaller than either individual resistance,

$$\frac{1}{R} = \frac{1}{R_1} + \frac{1}{R_2},$$

because the current has two different paths open to it.

An important example of resistances in parallel and series is the use of voltmeters and ammeters in resistive circuits. A voltmeter acts as a large resistance in parallel with the resistor across which the voltage drop is being measured. The fact that its resistance is not infinite means that it alters the circuit it is being used to investigate, producing a lower equivalent resistance. An ammeter acts as a small resistance in series with the circuit through which the current is to be determined. Its resistance is not quite zero, which leads to an increase in the resistance of the circuit being tested.

Chapter 10, Fields, page 579

Newton conceived of a universe where forces reached across space instantaneously, but we now

know that there is a delay in time before a change in the configuration of mass and charge in one corner of the universe will make itself felt as a change in the forces experienced far away. We imagine the outward spread of such a change as a ripple in an invisible universe-filling *field of force*.

As an alternative to our earlier energy-based definition, we can define the *gravitational field* at a given point as the force per unit mass exerted on objects inserted at that point, and likewise the *electric field* is defined as the force per unit charge. These fields are vectors, and the fields generated by multiple sources add according to the rules of vector addition.

The relationship between the electric field and the voltage is

$$\begin{aligned}\frac{\partial V}{\partial x} &= -E_x \\ \frac{\partial V}{\partial y} &= -E_y \\ \frac{\partial V}{\partial z} &= -E_z,\end{aligned}$$

which can be notated more compactly as a gradient,

$$\mathbf{E} = -\nabla V.$$

Fields of force contain energy, and the density of energy is proportional to the square of the magnitude of the field,

$$\begin{aligned}dU_g &= -\frac{1}{8\pi G}g^2 dv \\ dU_e &= \frac{1}{8\pi k}E^2 dv \\ dU_m &\propto B^2 dv\end{aligned}$$

The equation for the energy stored in the magnetic field is given explicitly in the next chapter; for now, we only need the fact that it behaves in the same general way as the first two equations: the energy density is proportional to the square of the field. In the case of static electric fields, we can calculate potential energy either using the previous definition in terms of mechanical work or by calculating the energy stored in the fields. If the fields are not static, the old method gives incorrect results and the new one must be used.

Capacitance, C , and inductance, L , are defined as

$$U_C = \frac{1}{2C}q^2$$

and

$$U_L = \frac{L}{2}I^2,$$

measured in units of farads and henries, respectively. The voltage across a capacitor or inductor is given by

$$V_C = \frac{q}{C}$$

or

$$|V_L| = \left| L \frac{dI}{dt} \right|.$$

In the equation for the inductor, the direction of the voltage drop (plus or minus sign) is such that the inductor resists the change in current. Although the equation for the voltage across an inductor follows directly from fundamental arguments concerning the energy stored in the magnetic field, the result is a surprise: the voltage drop implies the existence of electric fields which are not created by charges. This is an *induced electric field*, discussed in more detail in the next chapter.

A series LRC circuit exhibits *oscillation*, and, if driven by an external voltage, resonates. The Q of the circuit relates to the resistance value. For large Q , the resonant frequency is

$$\omega \approx \frac{1}{\sqrt{LC}}.$$

A series RC or RL circuit exhibits exponential *decay*,

$$q = q_0 \exp\left(-\frac{t}{RC}\right)$$

or

$$I = I_0 \exp\left(-\frac{R}{L}t\right),$$

and the quantity RC or L/R is known as the time constant.

When driven by a sinusoidal AC voltage with amplitude \tilde{V} , a capacitor, resistor, or inductor responds with a current having amplitude

$$\tilde{I} = \frac{\tilde{V}}{Z},$$

where the *impedance*, Z , is a frequency-dependent quantity having units of ohms. In a capacitor, the current has a phase that is 90° ahead of the voltage, while in an inductor the current is 90° behind. We can represent these phase relationships by defining the impedances as complex numbers:

$$\begin{aligned} Z_C &= -\frac{i}{\omega C} \\ Z_R &= R \\ Z_L &= i\omega L \end{aligned}$$

The arguments of the complex impedances are to be interpreted as phase relationships between the oscillating voltages and currents. The complex impedances defined in this way combine in series and parallel according to the same rules as resistances.

When a voltage source is driving a load through a transmission line, the maximum power is delivered to the load when the impedances of the line and the load are matched.

Gauss' law states that for any region of space, the flux through the surface,

$$\Phi = \sum \mathbf{E}_j \cdot \mathbf{A}_j,$$

is related by

$$\Phi = 4\pi k q_{in}$$

to the charge enclosed within the surface.

Chapter 11, Electromagnetism, page 675

Relativity implies that there must be an interaction between moving charges and other moving charges. This *magnetic* interaction is in addition to the usual electrical one. The magnetic field can be defined in terms of the magnetic force exerted on a test charge,

$$\mathbf{F} = q\mathbf{v} \times \mathbf{B},$$

or, alternatively, in terms of the torque on a magnetic test dipole,

$$|B| = \frac{\tau}{|\mathbf{m}_t| \sin \theta},$$

where θ is the angle between the dipole vector and the field. The magnetic dipole moment \mathbf{m} of a loop of current has magnitude $m = IA$, and is in the (right-handed) direction perpendicular to the loop.

The magnetic field has no sources or sinks. Gauss' law for magnetism is

$$\Phi_B = 0.$$

The external magnetic field of a long, straight wire is

$$B = \frac{2kI}{c^2 R},$$

forming a right-handed circular pattern around the wire.

The energy of the magnetic field is

$$dU_m = \frac{c^2}{8\pi k} B^2 dv.$$

The magnetic field resulting from a set of currents can be computed by finding a set of dipoles that combine to give those currents. The field of a dipole is

$$B_z = \frac{km}{c^2} (3 \cos^2 \theta - 1) r^{-3}$$

$$B_R = \frac{km}{c^2} (3 \sin \theta \cos \theta) r^{-3},$$

which reduces to $B_z = km/c^2 r^3$ in the plane perpendicular to the dipole moment. By constructing a current loop out of dipoles, one can prove the *Biot-Savart law*,

$$d\mathbf{B} = \frac{kI d\ell \times \mathbf{r}}{c^2 r^3},$$

which gives the field when we integrate over a closed current loop. All of this is valid only for static magnetic fields.

Ampère's law is another way of relating static magnetic fields to the static currents that created them, and it is more easily extended to nonstatic fields than is the Biot-Savart law. Ampère's law states that the *circulation* of the magnetic field,

$$\Gamma_B = \sum \mathbf{s}_j \cdot \mathbf{B}_j,$$

around the edge of a surface is related to the current I_{through} passing through the surface,

$$\Gamma = \frac{4\pi k}{c^2} I_{\text{through}}.$$

In the general nonstatic case, the fundamental laws of physics governing electric and magnetic fields are *Maxwell's equations*, which state that for any closed surface, the fluxes through the surface are

$$\begin{aligned}\Phi_E &= 4\pi k q_{\text{in}} && \text{and} \\ \Phi_B &= 0.\end{aligned}$$

For any surface that is not closed, the circulations around the edges of the surface are given by

$$\begin{aligned}\Gamma_E &= -\frac{\partial \Phi_B}{\partial t} && \text{and} \\ c^2 \Gamma_B &= \frac{\partial \Phi_E}{\partial t} + 4\pi k I_{\text{through}}.\end{aligned}$$

The most important result of Maxwell's equations is the existence of *electromagnetic waves* which propagate at the velocity of light — that's what light is. The waves are transverse, and the electric and magnetic fields are perpendicular to each other. There are no purely electric or purely magnetic waves; their amplitudes are always related to one another by $B = E/c$. They propagate in the right-handed direction given by the cross product $\mathbf{E} \times \mathbf{B}$, and carry momentum $p = U/c$.

A complete statement of Maxwell's equations in the presence of electric and magnetic materials is as follows:

$$\begin{aligned}\Phi_D &= q_{\text{free}} \\ \Phi_B &= 0 \\ \Gamma_E &= -\frac{d\Phi_B}{dt} \\ \Gamma_H &= \frac{d\Phi_D}{dt} + I_{\text{free}},\end{aligned}$$

where the auxiliary fields \mathbf{D} and \mathbf{H} are defined as

$$\begin{aligned}\mathbf{D} &= \epsilon \mathbf{E} && \text{and} \\ \mathbf{H} &= \frac{\mathbf{B}}{\mu},\end{aligned}$$

and ϵ and μ are the permittivity and permeability of the substance.

Chapter 12, Optics, page 765

The ray model of light: We can understand many phenomena involving light without having

to use sophisticated models such as the wave model or the particle model. Instead, we simply describe light according to the path it takes, which we call a ray. The ray model of light is useful when light is interacting with material objects that are much larger than a wavelength of light. Since a wavelength of visible light is so short compared to the human scale of existence, the ray model is useful in many practical cases.

A smooth surface produces specular reflection, in which the reflected ray exits at the same angle with respect to the normal as that of the incoming ray. A rough surface gives diffuse reflection, where a single ray of light is divided up into many weaker reflected rays going in many directions.

Images: A large class of optical devices, including lenses and flat and curved mirrors, operates by bending light rays to form an image. A real image is one for which the rays actually cross at each point of the image. A virtual image, such as the one formed behind a flat mirror, is one for which the rays only appear to have crossed at a point on the image. A real image can be projected onto a screen; a virtual one cannot.

Mirrors and lenses will generally make an image that is either smaller than or larger than the original object. The scaling factor is called the magnification. In many situations, the angular magnification is more important than the actual magnification.

Every lens or mirror has a property called the focal length, which is defined as the distance from the lens or mirror to the image it forms of an object that is infinitely far away. A stronger lens or mirror has a shorter focal length.

Locating images: The relationship between the locations of an object and its image formed by a lens or mirror can always be expressed by equations of the form

$$\begin{aligned}\theta_f &= \pm\theta_i \pm \theta_o \\ \frac{1}{f} &= \pm\frac{1}{d_i} \pm \frac{1}{d_o}.\end{aligned}$$

The choice of plus and minus signs depends on whether we are dealing with a lens or a mirror, whether the lens or mirror is converging or diverging, and whether the image is real or virtual. A method for determining the plus and minus signs is as follows:

1. Use ray diagrams to decide whether θ_i and θ_o vary in the same way or in opposite ways. Based on this, decide whether the two signs in the equation are the same or opposite. If the signs are opposite, go on to step 2 to determine which is positive and which is negative.
2. If the signs are opposite, we need to decide which is the positive one and which is the negative. Since the focal angle is never negative, the smaller angle must be the one with a minus sign.

Once the correct form of the equation has been determined, the magnification can be found via the equation

$$M = \frac{d_i}{d_o}.$$

This equation expresses the idea that the entire image-world is shrunk consistently in all three dimensions.

Refraction: Refraction is a change in direction that occurs when a wave encounters the interface between two media. Together, refraction and reflection account for the basic principles behind nearly all optical devices.

Snell discovered the equation for refraction,

$$n_1 \sin \theta_1 = n_2 \sin \theta_2,$$

[angles measured with respect to the normal]

through experiments with light rays, long before light was proven to be a wave. Snell's law can be proven based on the geometrical behavior of waves. Here n is the index of refraction. Snell invented this quantity to describe the refractive properties of various substances, but it was later found to be related to the speed of light in the substance,

$$n = \frac{c}{v},$$

where c is the speed of light in a vacuum. In general a material's index of refraction is different for different wavelengths of light. Total internal reflection occurs when there is no angle that satisfies Snell's law.

Wave optics: Wave optics is a more general theory of light than ray optics. When light interacts with material objects that are much larger than one wavelength of the light, the ray model of light is approximately correct, but in other cases the wave model is required.

Huygens' principle states that, given a wavefront at one moment in time, the future behavior of the wave can be found by breaking the wavefront up into a large number of small, side-by-side wave peaks, each of which then creates a pattern of circular or spherical ripples. As these sets of ripples add together, the wave evolves and moves through space. Since Huygens' principle is a purely geometrical construction, diffraction effects obey a simple scaling rule: the behavior is unchanged if the wavelength and the dimensions of the diffracting objects are both scaled up or down by the same factor. If we wish to predict the angles at which various features of the diffraction pattern radiate out, scaling requires that these angles depend only on the unitless ratio λ/d , where d is the size of some feature of the diffracting object.

Double-slit diffraction is easily analyzed using Huygens' principle if the slits are narrower than one wavelength. We need only construct two sets of ripples, one spreading out from each slit. The angles of the maxima (brightest points in the bright fringes) and minima (darkest points in the dark fringes) are given by the equation

$$\frac{\lambda}{d} = \frac{\sin \theta}{m},$$

where d is the center-to-center spacing of the slits, and m is an integer at a maximum or an integer plus $1/2$ at a minimum.

If some feature of a diffracting object is repeated, the diffraction fringes remain in the same places, but become narrower with each repetition. By repeating a double-slit pattern hundreds or thousands of times, we obtain a diffraction grating.

A single slit can produce diffraction fringes if it is larger than one wavelength. Many practical instances of diffraction can be interpreted as single-slit diffraction, e.g., diffraction in telescopes. The main thing to realize about single-slit diffraction is that it exhibits the same kind of relationship between λ , d , and angles of fringes as in any other type of diffraction.

Chapter 13, Quantum Physics, page 857

Quantum physics differs from classical physics in many ways, the most dramatic of which is that certain processes at the atomic level, such as radioactive decay, are random rather than deterministic. There is a method to the madness, however: quantum physics still rules out any process that violates conservation laws, and it also offers methods for calculating probabilities numerically. The most important of these generic methods is the law of independent probabilities, which states that if two random events are not related in any way, then the probability that they will both occur equals the product of the two probabilities,

$$\begin{aligned} \text{probability of A and B} \\ = & P_A P_B \quad [\text{if A and B are independent}]. \end{aligned}$$

When discussing a random variable x that can take on a continuous range of values, we cannot assign any finite probability to any particular value. Instead, we define the probability distribution $D(x)$, defined so that its integral over some range of x gives the probability of that range.

In radioactive decay, the time that a radioactive atom has a 50% chance of surviving is called the half-life, $t_{1/2}$. The probability of surviving for two half-lives is $(1/2)(1/2) = 1/4$, and so on. In general, the probability of surviving a time t is given by

$$P_{\text{surv}}(t) = 0.5^{t/t_{1/2}}.$$

Related quantities such as the rate of decay and probability distribution for the time of decay are given by the same type of exponential function, but multiplied by certain constant factors.

Around the turn of the twentieth century, experiments began to show problems with the classical wave theory of light. In any experiment sensitive enough to detect very small amounts of light energy, it becomes clear that light energy cannot be divided into chunks smaller than a certain amount. Measurements involving the photoelectric effect demonstrate that this smallest unit of light energy equals hf , where f is the frequency of the light and h is a number known as Planck's constant. We say that light energy is quantized in units of hf , and we interpret this quantization as evidence that light has particle properties as well as wave properties. Particles of light are called photons.

The only method of reconciling the wave and particle natures of light that has stood the test of experiment is the probability interpretation: the probability that the particle is at a given location is proportional to the square of the amplitude of the wave at that location.

One important consequence of wave-particle duality is that we must abandon the concept of the path the particle takes through space. To hold on to this concept, we would have to contradict the well established wave nature of light, since a wave can spread out in every direction simultaneously.

Light is both a particle and a wave. Matter is both a particle and a wave. The equations that connect the particle and wave properties are the same in all cases:

$$\begin{aligned} E &= hf \\ p &= h/\lambda \end{aligned}$$

Unlike the electric and magnetic fields that make up a photon-wave, the electron wavefunction is not directly measurable. Only the square of the wavefunction, which relates to probability, has direct physical significance.

A particle that is bound within a certain region of space is a standing wave in terms of quantum physics. The two equations above can then be applied to the standing wave to yield some important general observations about bound particles:

1. The particle's energy is quantized (can only have certain values).
2. The particle has a minimum energy.
3. The smaller the space in which the particle is confined, the higher its kinetic energy must be.

These immediately resolve the difficulties that classical physics had encountered in explaining observations such as the discrete spectra of atoms, the fact that atoms don't collapse by radiating away their energy, and the formation of chemical bonds.

A standing wave confined to a small space must have a short wavelength, which corresponds to a large momentum in quantum physics. Since a standing wave consists of a superposition of two traveling waves moving in opposite directions, this large momentum should actually be interpreted as an equal mixture of two possible momenta: a large momentum to the left, or a large momentum to the right. Thus it is not possible for a quantum wave-particle to be confined to a small space without making its momentum very uncertain. In general, the Heisenberg uncertainty principle states that it is not possible to know the position and momentum of a particle simultaneously with perfect accuracy. The uncertainties in these two quantities must satisfy the approximate inequality

$$\Delta p \Delta x \gtrsim h.$$

When an electron is subjected to electric forces, its wavelength cannot be constant. The "wavelength" to be used in the equation $p = h/\lambda$ should be thought of as the wavelength of the sine wave that most closely approximates the curvature of the wavefunction at a specific point.

Infinite curvature is not physically possible, so realistic wavefunctions cannot have kinks in them, and cannot just cut off abruptly at the edge of a region where the particle's energy would be insufficient to penetrate according to classical physics. Instead, the wavefunction "tails off" in the classically forbidden region, and as a consequence it is possible for particles to "tunnel" through regions where according to classical physics they should not be able to penetrate. If this quantum tunneling effect did not exist, there would be no fusion reactions to power our sun, because the energies of the nuclei would be insufficient to overcome the electrical repulsion between them.

Hydrogen, with one proton and one electron, is the simplest atom, and more complex atoms can often be analyzed to a reasonably good approximation by assuming their electrons occupy states that have the same structure as the hydrogen atom's. The electron in a hydrogen atom exchanges very little energy or angular momentum with the proton, so its energy and angular momentum are nearly constant, and can be used to classify its states. The energy of a hydrogen state depends only on its n quantum number.

In quantum physics, the angular momentum of a particle moving in a plane is quantized in units of \hbar . Atoms are three-dimensional, however, so the question naturally arises of how to deal with angular momentum in three dimensions. In three dimensions, angular momentum is a vector in the direction perpendicular to the plane of motion, such that the motion appears clockwise if viewed along the direction of the vector. Since angular momentum depends on

both position and momentum, the Heisenberg uncertainty principle limits the accuracy with which one can know it. The most that can be known about an angular momentum vector is its magnitude and one of its three vector components, both of which are quantized in units of \hbar .

In addition to the angular momentum that an electron carries by virtue of its motion through space, it possesses an intrinsic angular momentum with a magnitude of $\hbar/2$. Protons and neutrons also have spins of $\hbar/2$, while the photon has a spin equal to \hbar .

Particles with half-integer spin obey the Pauli exclusion principle: only one such particle can exist in a given state, i.e., with a given combination of quantum numbers.

We can enumerate the lowest-energy states of hydrogen as follows:

$n = 1, \quad \ell = 0, \quad \ell_z = 0,$	$s_z = +1/2 \text{ or } -1/2$	two states
$n = 2, \quad \ell = 0, \quad \ell_z = 0,$	$s_z = +1/2 \text{ or } -1/2$	two states
$n = 2, \quad \ell = 1, \quad \ell_z = -1, 0, \text{ or } 1, \quad s_z = +1/2 \text{ or } -1/2$		six states
...		...

The periodic table can be understood in terms of the filling of these states. The nonreactive noble gases are those atoms in which the electrons are exactly sufficient to fill all the states up to a given n value. The most reactive elements are those with one more electron than a noble gas element, which can release a great deal of energy by giving away their high-energy electron, and those with one electron fewer than a noble gas, which release energy by accepting an electron.

Index

- aberration, 798
 - chromatic, 811
- absolute zero, 75
- absorption, 769, 900
 - of waves, 374
- acceleration, 67
- adiabatic, 343
- alchemy, 15, 474
- alpha decay, 511
 - nature of emitted particle, 497
- alpha particle, *see* alpha decay
- ammeter, 536
- Ampère's law, 702
- ampere (unit), 533
- amplitude, 116
- analytic addition of vectors, 203
- anamorph, 844
- angular acceleration, 272
- angular frequency, 118
- angular magnification, 786
- angular momentum, 253
 - and the uncertainty principle, 923
 - in three dimensions, 923
 - of a particle in two dimensions, 254
 - of light, 254, 735
 - quantization of, 921
- angular velocity, 271
- antielectron, 513
- antimatter, 513
- Archimedean spiral, 760
- Archimedes' principle, 85, 207
- area
 - operational definition, 34
 - scaling of, 35
- Aristotle, 192
- asteroid, 193
- astrology, 15
- atom, 480
 - raisin-cookie model of, 492
- atomic number
 - defined, 501
- Atomism, 480
- atoms, 66, 68
 - helium, 939
 - lithium, 940
 - sodium, 941
 - with many electrons, 939
- averages, 861
 - rule for calculating, 861
- Avogadro's number, 319
- Bacon, Francis, 19
- ballast, 623
- Balmer, Johann, 930
- basis for a vector space, 968
- basis of a vector space, 968
- Bernoulli, Johann, 93
- beta decay, 512
 - nature of emitted particle, 497
- beta particle, *see* beta decay
- Big Bang, 69, 108
 - and the arrow of time, 337
 - described in general relativity, 453
 - evidence for, 370
- binding energy
 - nuclear, 517
- Biot-Savart law, 700
- black hole, 325, 434, 450
 - event horizon, 450
 - formation, 452
 - information paradox, 451
 - singularity, 452
- Bohr
 - Niels, 816
- Bohr, Niels, 931
- Boltzmann's constant, 319
- bond, *see* chemical bonds
- bottomonium, 951
- bound states, 899
- box
 - particle in a, 899
- brachistochrone, 94
- brightness of light, 771
- Brownian motion, 484
- buoyancy, 85, 207
- calorie
 - unit, 80

capacitor, 613
 capacitance, 613
 spherical, 609
carbon-14 dating, 867
Carnot engine, 321
cathode rays, 17, 488
causality, 398
Celsius (unit), 314
Celsius scale, 74
center of mass, 142
center of mass frame, 147
centi- (metric prefix), 24
cgs units, 1030
Chadwick, James, 140
chain reaction, 511
charge, 475
 conservation of, 477
 quantization of, 485
charmonium, 951
chemical bonds
 quantum explanation for hydrogen, 901
chemical reactions, 68
Chernobyl, 519
choice of axis theorem, 258
 proof, 1028
circuit, 535
 complete, 535
 open, 536
 parallel, 549
 series, 549
 short, 547
circular motion, 213
 inward force, 214
 no forward force, 214
 no outward force, 214
circular orbit, 98
circular orbit in a magnetic field, 684
classical physics, 858
climate change, 114, 522
closed system, 58
CMB, 109
coefficient of kinetic friction, 157
coefficient of static friction, 157
collision
 elastic, 139
 inelastic, 139
 totally inelastic, 148
color, 807
comet, 135
complete circuit, 535
complex numbers, 627
 in quantum physics, 913
component, 192
Compton scattering, 462
conduction of heat, 113
conductivity, 737
conductor
 defined, 542
conservation law, 56
conservation of mass, 480
convection, 113
converging, 783
conversions of units, 28
Copenhagen approximation, 889
 nonunitary, 973
Copenhagen interpretation, *see* Copenhagen approximation
correspondence principle, 69, 816
 defined, 397
 for mass-energy, 437
 for relativistic momentum, 432
 for time dilation, 397
cosmic censorship, 454
cosmic microwave background, 109, 455
cosmological constant, 108
coulomb (unit), 476
Coulomb's law, 477
coupling constant, 679
critically damped, 190
Crookes, William, 483
cross product, 286
 uniqueness, 1027
current
 defined, 533
current density, 612, 707
curved spacetime, 445
cyclotron, 749
 cyclotron frequency, 749
damped oscillations, 176
 critically damped, 190
 overdamped
 electrical, 624
 mechanical, 189
damping
 critical, 190

dark energy, 108, 456
 dark matter, 456
 Darwin, 18
 Darwin, Charles, 858
 Davisson
 C.J., 892
 de Broglie
 Louis, 892
 decay
 exponential, 865
 decoherence, 887, 1002
 Deep Space 1, 134
 definitions
 conceptual, 57
 operational, 57
 degeneracy, 922
 delta function (Dirac), 981
 derivative
 partial, 220
 Descartes, René, 131, 191
 Dialogues Concerning the Two New Sciences,
 36
 diamagnetism, 742
 differential mode, 740
 diffraction
 defined, 814
 double-slit, 818
 fringe, 815
 scaling of, 816
 single-slit, 823
 diffraction grating, 823
 diffuse reflection, 770
 digital camera, 873
 dimension of a vector space, 968
 diopter, 793
 dipole
 electric, 586
 energy due to orientation, 590
 field of, 695
 magnetic, 680
 field of, 696
 dipole moment, 587
 Dirac delta function, 981
 dispersion, 811, 896
 dispersive waves, 367
 dissonance, 386
 divergence, 654
 DNA, 519
 Doppler effect
 in relativity, 373
 Doppler shift, 368
 for light, 426
 gravitational, 447
 dot product, 216
 relativistic, 425
 double-slit diffraction, 818
 duality
 wave-particle, 880
 Dulong-Petit law, 334, 352
 dyne (unit), 1030
 Eötvös, Roland, 61
 Eddington
 Arthur, 858
 Einstein
 and randomness, 858
 Einstein's ring, 444
 Einstein, Albert, 857, 872
 and Brownian motion, 484
 electric current
 defined, 533
 electric dipole, 586
 field of, 695
 electric field, 585
 energy density of, 607
 related to voltage, 592
 electric forces, 475
 electrolytes, 551
 electromagnetic wave, 726
 momentum of, 448
 electron, 491
 as a wave, 892
 spin of, 936
 wavefunction, 895
 electron capture, 512
 electron decay, 512
 electrostatic unit, 1030
 elements, chemical, 482
 emf, 715
 emission spectrum, 900
 Empedocles of Acragas, 766
 endoscope, 808
 energy, 73
 “free”, 73
 distinguished from force, 154
 equivalence to mass, 433

heat, 74
 kinetic, 76
 light, 74
 quantization of for bound states, 900
 energy density
 of electric field, 607
 of gravitational field, 611
 of magnetic field, 612, 693
 energy-momentum four vector, 437
 engine
 automobile, 340
 Carnot, 321
 heat, 307
 Otto cycle, 340
 Enlightenment, 858
 entanglement, 885, 1007
 of macroscopic objects, 888
 entropy
 macroscopic definition, 323
 microscopic definition, 329
 equilibrium, 86
 metastable, 87
 neutral, 86
 redefined, 265
 stable, 86
 unstable, 87
 equipartition theorem, 334
 equivalence principle, 446
 equivalent resistance
 of resistors in parallel, 556
 erg (unit), 78, 1030
 escape velocity, 102
 esu (electrostatic unit), 1030
 ether, 413
 Euclidean geometry, 443
 Euler's formula, 629
 Euler, Leonhard, 629
 event horizon, 450
 evolution, 801
 randomness in, 858
 exclusion principle, 940, 987
 exponential decay, 865
 rate of, 868
 eye
 evolution of, 801
 human, 802
 farad
 defined, 614
 Faraday, Michael, 531
 types of electricity, 532
 Fermat's principle, *see* least time, principle of
 ferrite bead, 741
 ferromagnetism, 744
 field
 electric, 585
 gravitational, 581, 582
 fields
 superposition of, 583
 fields of force, 579
 flatworm, 801
 fluid
 defined, 208
 fluorescent light, 623
 flux
 additivity by charge, 645
 additivity by region, 645
 defined, 642
 in Gauss' theorem, 644
 focal angle, 791
 focal length, 792
 focal point, 792
 force
 analysis of forces, 160
 defined, 149
 distinguished from energy, 154
 fields of, 579
 normal, 156
 transmission, 162
 Foucault, 65
 four-vector, 425
 energy-momentum, 437
 Fourier's theorem, 368
 fourier-spectra, 386
 frame of reference, 63
 inertial, 63
 in general relativity, 449
 Franklin, Benjamin
 definition of signs of charge, 477
 French Revolution, 24
 frequency, 118
 of waves, 365
 friction
 fluid, 159
 kinetic, 156
 static, 156

fringe
 diffraction, 815
 full width at half maximum, 865
 full width at half-maximum, 184
 fundamental, 386
 fundamental theorem of algebra, 628
 FWHM, 184, 865

 Galileo, 767
 Galilean relativity, 62, 195
 inertial and gravitational mass, 61
 Galileo Galilei, 36
 gamma decay
 nature of emitted particle, 497
 gamma ray, *see* gamma decay
 pair production, 439
 garage paradox, 410
 gas
 spectrum of, 900
 gas discharge tube, 623
 gauss (unit), 1030
 Gauss' law, 650
 differential form, 654
 Gauss' theorem, 644
 for gravity, 650
 proof of, 649
 Gaussian pillbox, 651
 general relativity, 443
 generator, 622, 716
 geothermal vents, 857
 Germer, L., 892
 GFI, 689
 Gisin's theorem, 918
 global warming, 114, 522
 goiters, 865
 gradient, 220
 graphical addition of vectors, 203
 gravitational field, 82, 581, 582
 energy density of, 611
 gravitational time dilation, 447
 gravitational waves, 584
 Gravity Probe B, 444
 ground fault interrupter, 689
 group velocity, 898

 half-life, 865
 Halley's comet, 135
 Hamiltonian, 992
 handedness, 687

 harmonics, 386
 Hawking radiation, 325
 Hawking singularity theorem, 454
 Hawking, Stephen, 454
 heat, 74, 308, 315
 compared to temperature, 74
 compared to thermal energy, 75
 heat capacity
 at constant pressure, 342
 at constant volume, 342
 heat engine, 307
 heat transfer
 by conduction, 113
 by convection, 113
 by radiation, 113
 Heisenberg
 Werner, 902
 Heisenberg uncertainty principle, 902
 in three dimensions, 923
 helium, 939
 Helmholtz resonator, 344
 hermitian operator, 986
 Hertz
 Heinrich, 732
 Hertz, Heinrich, 875
 Heinrich, 818
 Hiroshima, 520
 homogeneity of spacetime, 403
 Hooke, 474
 Hooke's law, 173
 hormesis, 521
 Hubble, Edwin, 370
 Hugo, Victor, 473
 Huygens' principle, 817
 hydrogen atom, 927
 angular momentum in, 921
 classification of states, 920
 energies of states in, 929
 energy in, 921
 momentum in, 921
 quantum numbers, 927
 hydrogen molecule, *see* chemical bonds
 hysteresis, 745

 ideal gas law, 319
 images
 formed by curved mirrors, 783
 formed by plane mirrors, 780

- location of, 790
- method of (electrostatics), 657
- of images, 785
- real, 784
- virtual, 780
- impedance, 630
 - of an inductor, 633
- impedance matching, 637, 741
- incoherent light, 815
- independence
 - statistical, 859, 860
- independent probabilities
 - law of, 860
- index of refraction
 - defined, 804
 - related to speed of light, 805
- inductance
 - defined, 615
- induction, 622
- inductor, 613
 - inductance, 613
- inertial frame of reference, 63
- information paradox, 451
- inner product, 425
 - in quantum mechanics, 984
 - on a general vector space, 968
- insulator
 - defined, 542
- invariance
 - rotational, 195
- inverted reflection, 376
- Io, 768
- iodine, 865
- ion drive, 134
- isotopes, 508
- Ives-Stilwell experiment, 428
- Jeans
 - James, 858
- joule (unit), 74
- Joule, James, 73
 - paddlewheel experiment, 76
- junction rule, 555
- Jupiter, 768
- kelvin (unit), 314
- Kelvin scale, 75
- Kepler's laws, 96
- Keynes, John Maynard, 474
- kilo- (metric prefix), 24
- kilogram, 25, 56
 - standard, 57
- kinetic energy, 76
 - compared to momentum, 136
- kinetic energy theorem, 169
- kinetic friction, 156
 - coefficient of, 157
- Lagrange, Joseph-Louis, 55
- Laplace, 16
- Laplace, Pierre Simon de, 857
- Laplacian, 657, 912
- Lavoisier, Pierre-André
 - conservation of mass, 59
 - execution, 55
- least time, principle of, 778, 811, 827
- Leibniz, 94
- lens, 808
- lensmaker's equation, 810
- light, 16
 - absorption of, 769
 - angular momentum of, 735
 - brightness of, 771
 - defined, 482
 - Doppler shift for, 426
 - electromagnetic wave, 726
 - momentum of, 135, 448, 460, 732
 - particle model of, 771
 - ray model of, 771
 - speed of, 767
 - wave model of, 771
 - waves, 364
- light cone, 420
- lightlike, 420
- LIGO, 585
- line integral, 220
- linear independence, 968
- linear no-threshold, 521
- linear operator, 968
- Lipkin linkage, 843
- LNT, 521
- loop rule, 560
- Lorentz invariance, 423
- Lorentz transformation, 404
- Lorentz, Hendrik, 404
- LRC circuit, 639
- lumped-circuit approximation

for capacitors, 603

magnetic dipole, 680
field of, 696

magnetic field
defined, 680, 681
energy density of, 612, 693
long, straight wire
found using Biot-Savart law, 752

magnetic monopoles, 685

magnification
angular, 786
by a converging mirror, 783

many-worlds approximation, 889

many-worlds interpretation, *see* many-worlds approximation

mass
conservation of, 56, 480
equivalence to energy, 433
gravitational, 57
inertial, 57
quantization of, 69

mass-energy
conservation of, 435
correspondence principle, 437
of a moving particle, 436

matter, 16
as a wave, 891
defined, 482

Maxwell's equations, 724
for static fields, 705
in cgs units, 1030
in differential form, 1029

Maxwell, James Clerk, 818

mechanical system, 132

median, 865

mega- (metric prefix), 24

Mendeleev, Dmitri, 483

meter (unit), 56

metric system, 24, 56
prefixes, 24, 1072

Michelson-Morley experiment, 413

micro- (metric prefix), 24

microwaves, 16

milli- (metric prefix), 24

Millikan, Robert, 485

mirror
converging, 790

mks units, 25

molecules
nonexistence in classical physics, 891

mollusc, 802

moment
dipole, 587

moment of inertia, 274
tabulated for various shapes, 281

momentum, 132, 133
compared to kinetic energy, 136
nonmechanical, 135
of light, 135, 448, 460, 732
relativistic, 429, 437

monopoles
magnetic, 685

MRI (magnetic resonance imaging), *see* NMR,
see NMR

naked singularity, 454

nano- (metric prefix), 24

natural units, 406

nautilus, 801

neutral (electrically), 477

neutral equilibrium, 86

neutron
discovery of, 140
spin of, 936

Newton
Isaac, 857

newton (unit), 149

Newton, Isaac, 94
alchemy, 473
definition of time, 26
Newtonian telescope, 785
particle theory of light, 817

Nichols-Hull experiment on momentum of light, 734

NMR (nuclear magnetic resonance), 684, 994

no-cloning theorem, 1005, 1009

normal force, 156

normal operator, 986

normalization, 861

nuclear forces, 509, 688

nuclear reactions, 68

nucleus
discovery, 499

Ohm's law, 542

ohmic

defined, 542
op-amp, 619
open circuit, 536
operational amplifier (op-amp), 619
operational definitions, 24
orbit
 circular, 98
order-of-magnitude estimates, 43
oscillations, 115
 damped, 176
 overdamped
 electrical, 624
 mechanical, 189
 steady state, 180
Otto cycle, 340
Otto, Nikolaus, 340
overdamped oscillations
 electrical, 624
 mechanical, 189
ozone layer, 872

paddlewheel experiment, 76
pair production, 439, 462
parallel axis theorem, 276, 302
parallel circuit
 defined, 549
paramagnetism, 742
parity
 operator in quantum mechanics, 983
Parmenides, 55
partial derivative, 220, 654
particle
 definition of, 880
particle in a box, 899
particle model of light, 771, 818
pascal
 unit, 309
pascal (unit), 208
path of a photon undefined, 881
Pauli exclusion principle, 18, 940, 987
Peaucellier linkage, 843
Pelton waterwheel, 193
Penrose singularity theorem, 453
Penrose, Roger, 453
period
 of waves, 365
periodic table, 483, 501, 940
permeability, 739

permittivity, 738
perpetual motion machine, 73
phase in quantum mechanics
 not observable, 896, 917, 969, 978, 982
phase velocity, 898
photoelectric effect, 875
photon
 Einstein's early theory, 874
 energy of, 877
 in three dimensions, 889
 spin of, 936
physics, 16
pillbox
 Gaussian, 651
pilot wave theory, 970
Planck's constant, 877
Planck, Max, 877
polarization, 728
Pope, 36
positron, 435, 513
positron decay, 512
potential
 electrical, 538
Pound-Rebka experiment, 448
power, 80
 electrical, 538
Poynting vector, 756
praxinoscope, 781
pressure, 308
 as a function of depth, 208
 defined, 208
probabilities
 addition of, 861
 normalization of, 861
probability distributions
 averages of, 864
 widths of, 864
probability distributions, 863
probability interpretation, 881
protein molecules, 920
proton
 spin of, 936
Pythagoras, 766

quality factor, 180
quantization, 485
 of mass, 69
quantum dot, 899

quantum mechanics
 measurement problem, 1001

quantum moat, 921

quantum numbers, 928

quantum physics, 858

quark, 698, 951

radar, 872

radiation hormesis, 521

radiation of heat, 113

radio, 872

radio waves, 16

raisin cookie model, 492

randomness, 858

ray diagrams, 773

ray model of light, 771, 818

RC circuit, 624

RC time constant, 625

reactions
 chemical, 68
 nuclear, 68

reductionism, 19

reflection
 diffuse, 770
 of waves, 374
 specular, 774

reflections
 inverted and uninverted, 376

refraction
 and color, 807
 defined, 802

relativity
 Galilean, 62, 195
 general, 443

Renaissance, 13

repetition of diffracting objects, 822

resistance
 defined, 542
 in parallel, 555
 in series, 559

resistivity
 defined, 561

resistor, 547

resistors
 in parallel, 556

retina, 785

reversibility, 776

RHIC accelerator, 409

RL circuit, 625

RMS (root mean square), 636

Roemer, 768

root mean square, 636

rotational invariance, 195

Russell
 Bertrand, 858

Rutherford
 characterization of alpha particles, 497
 discovery of nucleus, 499
 Ernest, 857

scalar
 defined, 197

scaling, 35

schematic, 554

schematics, 554

Schrödinger equation, 906
 time-dependent, 970, 991

scientific method, 14

sea-of-arrows representation, 583

second (unit), 25, 56

separability
 of a quantum state, 1007

Schrödinger equation, 976

series circuit
 defined, 549

shell theorem, 102

short circuit
 defined, 547

SI units, 25, 78

Sievert (unit), 519

sigma notation, 143

significant figures, 30

simple machine, 172

single-slit
 diffraction, 823

singularity
 Big Bang, 453
 black hole, 452
 naked, 454

singularity theorem
 Hawking, 454
 Penrose, 453

sinks in fields, 583

Sirius, 900

skin depth, 737

Snell's law, 803

derivation of, 806
mechanical model of, 805
sodium, 941
solar constant, 756
solar sail, 208
solenoid, 614
 magnetic field of, 703
Sommerfeld, Arnold, 352
sound
 speed of, 389
 waves, 363
sources of fields, 583
spacelike, 420
spacetime
 curvature of, 445
spark plug, 625
specific heat, 74
 electrons' contribution, 352
 monoatomic ideal gas, 333
 solids, 334
spectrum
 absorption, 900
 emission, 900
spherical harmonics, 925
spin, 698, 936
 neutron's, 936
 of electron, 936
 photon's, 936
 proton's, 936
spin theorem, 259
 proof, 1029
spiral
 Archimedean, 760
spring constant, 115
Squid, 802
stability, 86
standing wave, 386
standing waves, 386
Star Trek, 901
states
 bound, 899
static friction, 156
 coefficient of, 157
statvold, 1030
steady state, 180
strong nuclear force, 509
strong nuclear force, 509
superposition
 of waves, 354
superposition of fields, 583
Swift, Jonathan, 35
symmetry, 687
system
 closed, 58
telescope, 785, 824
temperature, 74, 308
 absolute zero, 75, 314
 Celsius, 314
 compared to heat, 74
 Kelvin, 314
 macroscopic definition, 314
 microscopic definition, 330
tension, 163, 171
tensor product
 of quantum states, 1008
Tesla
 Nikola, 182
tesla (unit), 680
thermal energy
 compared to heat, 75
thermodynamics, 307
 first law of, 308, 339
 laws of
 summarized, 339
 second law of, 324, 338, 339
 third law of, 339
 zeroth law of, 313, 339
thermometer, 314
Thomson, J.J.
 cathode ray experiments, 489
time
 arrow of, 338
time constant
 RC, 625
time dilation
 gravitational, 447
time reversal, 776
timelike, 420
Tolman-Stewart experiment, 595
torque
 defined, 260
 related to force, 261, 290
total internal reflection, 808
transformer, 622, 716
transmission

of waves, 374
 transmission of forces, 162
 triangle inequality, 425
 tunneling, 905
 twin paradox, 425
 ultraviolet light, 872
 uncertainty principle, 902
 in three dimensions, 923
 unitary evolution of the wavefunction, 972, 987
 unitary operator, 987
 units
 natural relativistic, 406
 nonmetric, 1072
 units, conversion of, 28
 unstable equilibrium, 87
 vector
 addition, 198
 defined, 197
 division by a scalar, 198
 dot product, 216
 four-vector, 425
 magnitude of, 198
 multiplication by a scalar, 198
 subtraction, 198
 vector addition
 analytic, 203
 graphical, 203
 vector cross product, 286
 vector product, cross, 286
 vector space, 967
 velocity
 addition of
 relativistic, 428, 461
 vector, 205
 group, 898
 phase, 898
 velocity filter, 684
 vision, 766
 volt (unit)
 defined, 537
 voltage
 defined, 538
 related to electric field, 592
 volume
 operational definition, 34
 scaling of, 35
 Voyager space probe, 458
 water
 specific heat, 74
 wave
 definition of, 880
 dispersive, 811, 896
 electromagnetic, 726
 energy related to amplitude, 377
 light, 726
 wave model of light, 771, 818
 wave-particle duality, 880
 probability interpretation of, 881
 wavefunction
 complex numbers in, 913
 of the electron, 895
 wavelength, 366
 wavenumber, 998
 waves
 absorption of, 374
 dispersive, 367
 frequency of, 365
 gravitational, 584
 interference, 382
 light, 364
 medium not transported with, 356
 on a string, 360
 patterns, 358
 period of, 365
 reflection of, 374
 sound, 363
 standing, 386
 superposition of, 354
 transmission of, 374
 velocity of, 357
 wavelength, 366
 weak nuclear force, 512, 688
 whale songs, 376
 Wicked Witch of the West, 891
 Wigner, Eugene, 790
 work
 defined, 166
 work theorem, 168
 world-line, 425
 x-rays, 16
 Yarkovsky effect, 193
 Young, Thomas, 818
 zero-point motion, 965